



Model-free Prediction

Blink Sakulkeakulsuk

Four main Bellman Equations

Bellman Expectation Equation

⌋ value | Policy π

$$v_{\pi}(s) = \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s, A_t \sim \pi(s)]$$

$$q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$$

Bellman Optimality Equation

$$v^*(s) = \max_a \mathbb{E}[R_{t+1} + \gamma v^*(S_{t+1}) | S_t = s, A_t = a]$$

$$q^*(s, a) = \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} q^*(S_{t+1}, a') \mid S_t = s, A_t = a \right]$$

Prediction and Control

Prediction *

- A problem when we perform **policy evaluation**
- Estimating v_π or q_π *try to ค่า policy*

Control

- A problem when we perform **policy optimization**
- Estimating v^* or q^* *try to have process to optimize policy*

Evaluating Policy

$$\pi \geq \pi' \leftrightarrow v_{\pi}(s) \geq v_{\pi'}(s) , \forall s$$

One policy is better than or equal to another **if and only** if its value function is greater or equal to another.

Policy Evaluation

Given a policy, we want to estimate this

$$v_{\pi}(s) = \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | s, \pi]$$

We initialize all value to zero, and we iterate the equation above as an update

$$v_{k+1}(s) = \mathbb{E}[R_{t+1} + \gamma v_k(S_{t+1}) | s, \pi] \quad , \quad \forall s$$

Note that

$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_r \sum_{s'} p(r, s' | s, a) (r + \gamma v_{\pi}(s'))$$

If $v_{k+1}(s) = v_k(s), \forall s$, we solve v_{π} .

Policy Improvement

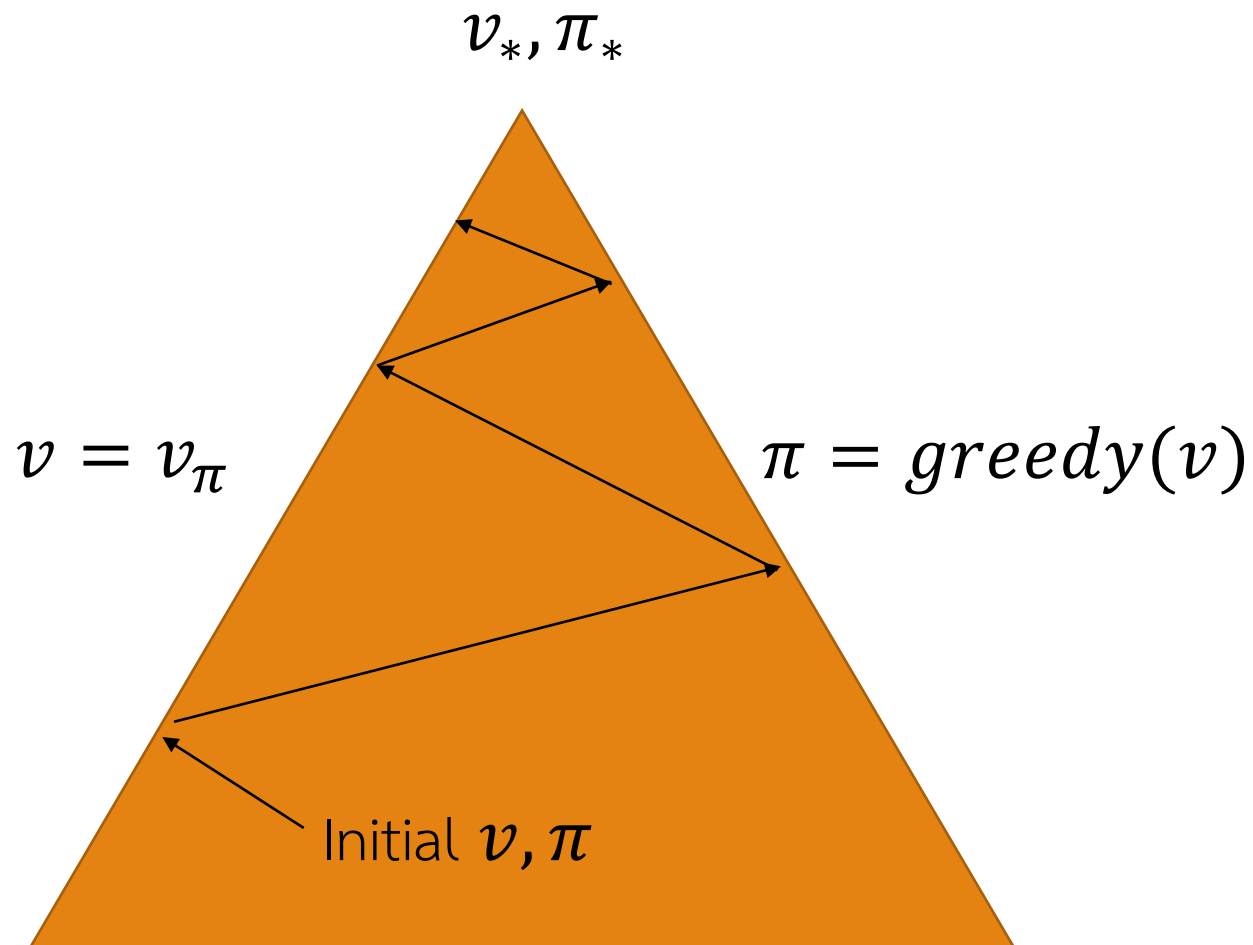
The example shows that we can improve a policy when we learn the value.

- In the example, we do not improve any policy

$$\begin{aligned}\forall s: \pi_{new}(s) &= \operatorname{argmax}_a q_{\pi}(s, a) \\ &= \operatorname{argmax}_a \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s, A_t = a]\end{aligned}$$

We use the new policy to evaluate, and repeat.

Policy Iteration



Starting with initial v, π

Loop do:

Perform policy evaluation,

$$v = v_\pi$$

Perform policy improvement,

$$\pi' \geq \pi$$

Value Iteration

Or we can learn the policy on the fly

- We update the policy every iteration

We can take the Bellman optimality equation as an update

$$v_{k+1}(s) \leftarrow \max_a \mathbb{E}[R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s, A_t = a] \quad , \quad \forall s$$

It is the same as **policy iteration** with $k = 1$. Basically, we improve the policy every iteration.

Model-Free Learning

Dynamic Model

What will happen if we don't know dynamic Model

$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_r \sum_{s'} p(r, s' | s, a) (r + \gamma v_{\pi}(s'))$$

The value function above contains the dynamic of the problem

- $p(r, s' | s, a)$ tells the transition model of the problem. This is essentially the dynamic of the problem.

ถ้าไม่รู้มันจะคำนวณทุกอย่างได้ไหม ?

Key Question: What if we do not know that? Can the agent still learn?

- Spoiler Alerted: Yes, it can. (Otherwise, we wouldn't have this class)

Our Plan

Last Lecture

Markov

- We solved a **known MDP** by dynamic programming

This Lecture

- **Model-free prediction** to **estimate values** in an **unknown MDP**
- A brief touch of Deep Reinforcement Learning

Next Lecture

- **Model-free control** to **optimize values** in an **unknown MDP**

Monte Carlo Algorithm

Monte Carlo Algorithm

เราให้ run 1 ครั้งจนจบ แล้วมาดูว่าเรายังรู้อะไรบ้าง

To learn without a model, we can **sample the experience**

- Let's just act until the end and learn from that

Monte Carlo is a sampling method

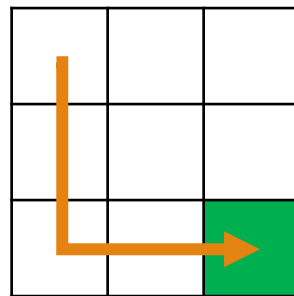
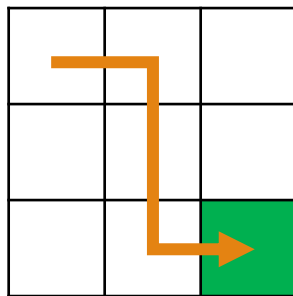
- Direct sampling of episodes
- Model-free
- No knowledge of MDP required

Monte Carlo Example

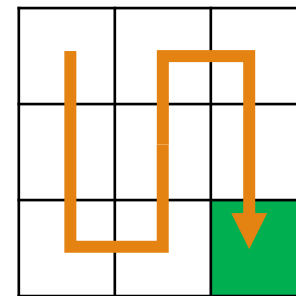
To sample the experience, given a policy π , the agent just plays out until termination.

For example, given a grid world and a robot that can move up, down, left, right, the samples are:

Episode 1 : run รอบ 1 Episode 2



Episode 3



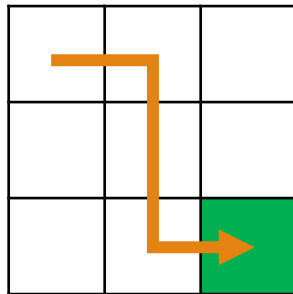
การที่ทำได้จนถึง goal state ไม่แปลว่า 3 Return

Monte Carlo Example

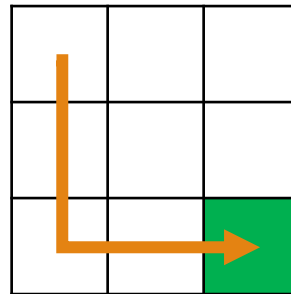
To sample the experience, given a policy π , the agent just plays out until termination.

For example, given a grid world and a robot that can move up, down, left, right, the samples are:

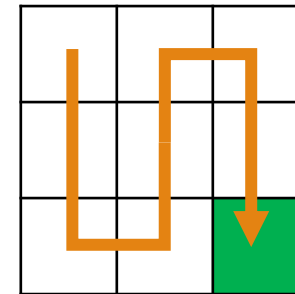
Episode 1



Episode 2



Episode 3



Monte-Carlo Policy Evaluation

Given the sequential decision problem, MC learns v_π from episodes under policy π

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

The return is calculated given an ending time T

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

↳ **return** คือผลรวมของ state value ไปถึง

The value function is then calculated from the expected return

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

We can use **sample average return** instead of **expected return**

- This is called **Monte Carlo policy evaluation**

กำหนด sat of γ
 block timestep เพื่อป้องกัน ∞ timestep แล้วไม่ถึง goal

Agent sampling ไปเรื่อย ๆ

โดย reward จะถูก avg เพื่อในอีกขั้น

True value

Part 1 Summary: Monte Carlo

ตรวจสอบเงื่อนไขว่าถึง goal

* เริ่มต้นจาก

การสุ่ม reward

Input: policy π , num_episodes

Initialize $N(s, a) = 0$ for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$ ตาม DL ที่วาง

Initialize $returns_sum(s, a) = 0$ for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

for $i \leftarrow 1$ to num_episodes do

Generate an episode $S_0, A_0, R_1, \dots, S_T$ using $\pi \rightarrow$ เก็บค่าทุก timestep เพื่อคำนวณ

for $t \leftarrow 0$ to $T - 1$ do

timestep

if (S_t, A_t) is a first visit (with return G_t) then

why first visit ?

$N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$ เพิ่ม $N(S_t, A_t)$

• เราจะไม่คำนวณจำนวนที่กลับมามา

$returns_sum(S_t, A_t) \leftarrow returns_sum(S_t, A_t) + G_t$ ซ่อมเดิม

$G_t + state 0 \neq G_t + state 1$

end

bc reward not equal except reward = 0

end Agent ตรวจสอบว่า target คืออะไร

value
or
ถ้าได้

$Q(s, a) \leftarrow returns_sum(s, a) / N(s, a)$ for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

คำนวณครั้งที่ถูก visit

return Q

* มีโอกาสเกิดวนลูป

Blackjack Example

200 States

6th 2nd HW 2

- Current hand (12-21)
- Dealer's hand (A-10)
- Our usable ace (yes, no)

Actions: stick, draw

Reward

- Stick: +1 if win, 0 if draw, -1 if lose
- Draw: -1 if lose, 0 otherwise

Transitions: draw if sum < 12

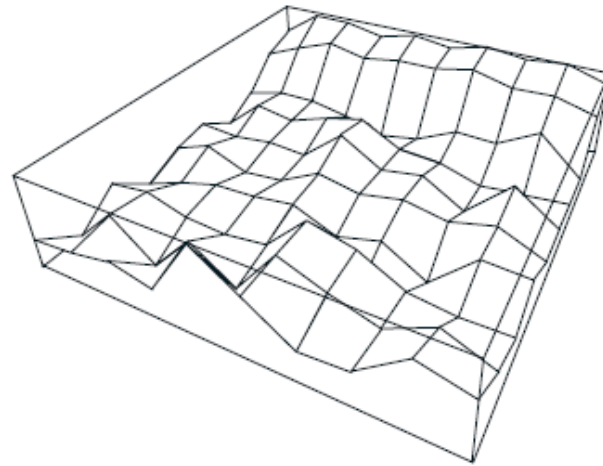
eng

Goal of HW 2

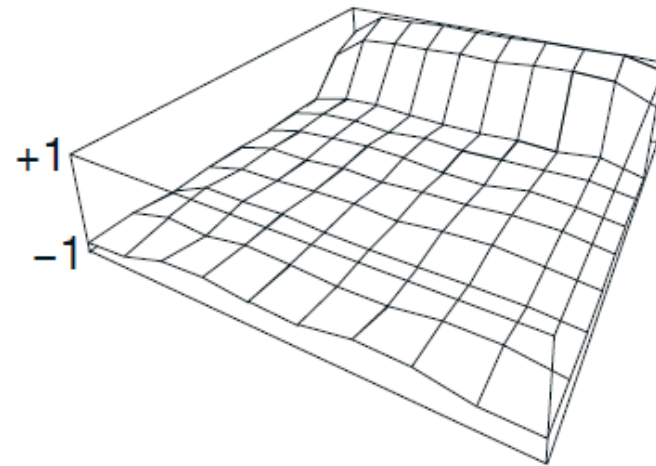
After 10,000 episodes

After 500,000 episodes

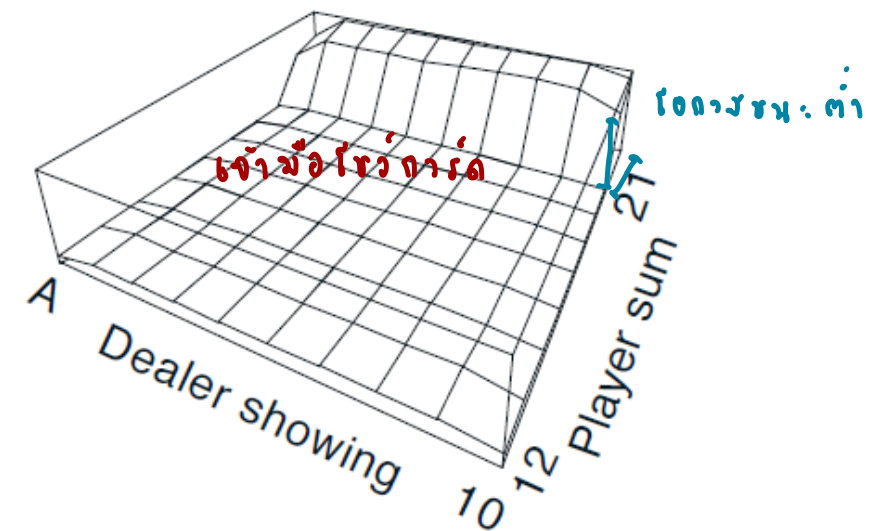
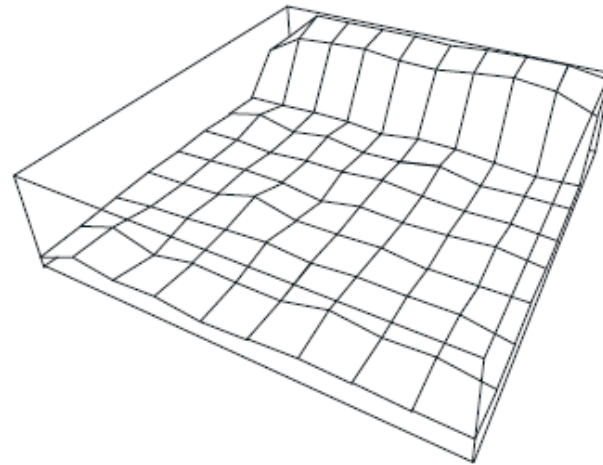
Usable
ace



+1
-1



No
usable
ace



Function Approximation

DL + RL

from Monte carlo

$Q(s,a)$ always add s,a

Value Function Approximation

In dynamic programming with MDP, we use **lookup tables** to learn value function.

- A mapping of state s to a value $v(s)$
- A mapping of state-action pair s, a to a state-action value $q(s, a)$

But we cannot do that for a large MDP

- Infeasible to store all information due to the **curse of dimensionality**
- **Too slow to learn** the value of each state **individually**
- States are often **not fully observable**

Value Function Approximation

For large MDP, we approximate the value function instead

$$\left. \begin{aligned} v_{\mathbf{w}}(s) &\approx v_{\pi}(s) \\ q_{\mathbf{w}}(s, a) &\approx q_{\pi}(s, a) \end{aligned} \right\} \begin{array}{l} \text{แทนที่: ค่าคาดหวัง} \\ \text{ของสถานะหรือสถานะ-การกระทำ} \end{array}$$

The agent will learn the value and update (learnable) parameter \mathbf{w} .^{vector}

- Monte Carlo Algorithm
- Temporal Difference Learning

เมื่อค่าที่ได้ (goal, return) มาเพื่อปรับใน \mathbf{w} ?

Agent State Update

For large MDP, if the environment state is not fully observable, we can update the agent state by

$$S_t = u_{\omega}(S_{t-1}, A_{t-1}, O_t)$$

with parameters $\omega \in \mathbb{R}^n$

Linear Value Function Approximation

Let's start with a simple, useful, and special case: a linear function

A state is represented by a **feature vector**

$$\mathbf{x}(s) = \begin{bmatrix} x_1(s) \\ \vdots \\ x_m(s) \end{bmatrix}$$

size m

$\mathbf{x} : S \Rightarrow \mathbb{R}^m$ is a mapping from an agent state to features

Linear Value Function Approximation

One way to approximate a value function is to find a linear combination of features.

$$v_{\mathbf{w}}(s) = \mathbf{w}^T \mathbf{x}(s) = \sum_{j=1}^n w_j x_j(s)$$

• η Model function:

linear regression

To learn the function, we can find an objective loss.

$$L(\mathbf{w}) = \mathbb{E}_{S \sim d} \left[\left(v_{\pi}(S) - \underbrace{\mathbf{w}^T \mathbf{x}(s)}_{\substack{\text{things we want} \\ \text{to learn}}} \right)^2 \right]$$

Linear Value Function Approximation

?

$$L(\mathbf{w}) = \mathbb{E}_{S \sim d} \left[\left(v_{\pi}(S) - \mathbf{w}^T \mathbf{x}(S) \right)^2 \right] \rightarrow \text{loss function of the return}$$

Then we can use stochastic gradient descent to update the parameter

- Stochastic gradient descent is a gradient descent algorithm that calculates loss of 1 input at a time
- Since we sample the problem, we can update the parameter using SGD

$$\nabla_{\mathbf{w}} v_{\mathbf{w}}(S_t) = \mathbf{x}(S_t) = \mathbf{x}_t \xrightarrow{\text{feature vector}} \Delta \mathbf{w} = \alpha \underbrace{\left(v_{\pi}(S_t) - v_{\mathbf{w}}(S_t) \right)}_{\text{step size}} \mathbf{x}_t$$

Update = step-size x prediction error x feature vector

Function Approximation Example

Let's consider a bandit problem. If we want to learn the value function, we can find it through the objective loss function.

- q is a parametric function with parameters \mathbf{w} . For example, q could be a neural network.

$$L(\mathbf{w}) = \frac{1}{2} \mathbb{E}[(R_{t+1} - q_{\mathbf{w}}(S_t, A_t))^2]$$

Handwritten notes:
 R_{t+1} is labeled v_{π}
 $q_{\mathbf{w}}(S_t, A_t)$ is labeled "parametric eq."
 $L(\mathbf{w})$ is labeled "loss" with a note "it's also" and an arrow pointing to the equation.

Then we update the weight by:

$$\begin{aligned} \mathbf{w}_{t+1} &= \mathbf{w}_t - \alpha \nabla_{\mathbf{w}_t} L(\mathbf{w}_t) \\ &= \mathbf{w}_t - \alpha \nabla_{\mathbf{w}_t} \frac{1}{2} \mathbb{E}[(R_{t+1} - q_{\mathbf{w}}(S_t, A_t))^2] \\ &= \mathbf{w}_t + \alpha \mathbb{E}[(R_{t+1} - q_{\mathbf{w}}(S_t, A_t)) \nabla_{\mathbf{w}_t} q_{\mathbf{w}}(S_t, A_t)] \end{aligned}$$

Handwritten notes:
 $\nabla_{\mathbf{w}_t}$ is labeled "หา gradient" (find gradient).

Function Approximation Example

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha \mathbb{E}[(R_{t+1} - q_{\mathbf{w}}(S_t, A_t)) \nabla_{\mathbf{w}_t} q_{\mathbf{w}}(S_t, A_t)]$$

If we use a linear function approximation,

$$q(s, a) = \mathbf{w}^T \mathbf{x}(s, a)$$

(Handwritten note: $\nabla_{\mathbf{w}_t}$)

Handwritten note: Prediction Problem in Model free env. so we sampling until get return, change

The gradient becomes

$$\nabla_{\mathbf{w}_t} q_{\mathbf{w}}(S_t, A_t) = \mathbf{x}(s, a)$$

Then the SGD update is

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha (R_{t+1} - q_{\mathbf{w}}(S_t, A_t)) \mathbf{x}(s, a)$$

Linear Update = step-size x prediction error x ~~feature vector~~
gradient

Non-linear Update = step-size x prediction error x ~~feature vector~~

Monte Carlo → Problem : so Agent จะทำอะไร , จะ run อย่างไร : ไม่รู้

MC algorithm can be used to learn value function

- However, the algorithm must wait until the end of an episode to learn something
- Return can have high variance

Alternatives?

- Yes! Temporal Difference Learning → รู้ก่อนจะจบ : จะทำอะไร , จะ run อย่างไร : ไม่รู้

Temporal Difference Learning

Temporal Difference Learning

Given a Bellman equation,

$$v_{\pi}(s) = \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s, A_t \sim \pi(S_t)]$$

We can update values by iterating using the following update function,

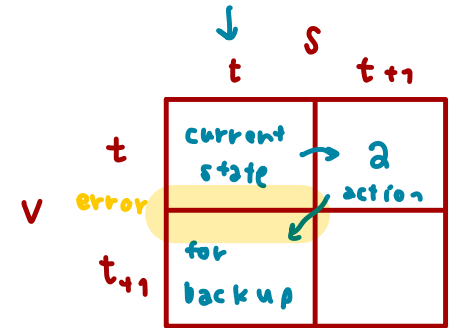
$$v_{k+1}(s) = \mathbb{E}[R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s, A_t \sim \pi(S_t)]$$

can change
form to
update value

Temporal Difference Learning

Monte Carlo - ทฤษฎีบทค่าเฉลี่ย

$$v_{k+1}(s) = \mathbb{E}[R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s, A_t \sim \pi(S_t)]$$



Instead of calculating the expected return, we can sample the value.

if we have noise and not using TDL, we avg out for a long time
if we have noise and not using TDL, the state = 10 for a long time

$$v_{t+1}(S_t) = R_{t+1} + \gamma v_t(S_{t+1})$$

we want to know
value of S_t after we action
 S_{th}

But instead of updating everything on the noisy value, we can update the value a little bit instead.

$$v_{t+1}(S_t) = v_t(S_t) + \underbrace{\alpha_t}_{\text{learning factor sth}} \underbrace{(R_{t+1} + \gamma v_t(S_{t+1}) - v_t(S_t))}_{\text{value that sampling - } v_t, \text{ an error between time step}}$$

MC Prediction vs TD Prediction

value function v_π of policy π

Prediction problem: learn v_π online from experience under policy π

Monte Carlo:

- Update value $v_n(S_t)$ with respect to sample return G_t

at episode end update

$$v_n(S_t) = v_n(S_t) + \alpha(G_t - v_n(S_t))$$

value

return

(return - value)

Temporal-difference learning:

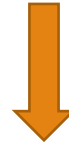
- Update value $v_t(S_t)$ with respect to estimated return $R_{t+1} + \gamma v(S_{t+1})$

$$v_{t+1}(S_t) = v_t(S_t) + \alpha_t(R_{t+1} + \gamma v_t(S_{t+1}) - v_t(S_t))$$

at each timestep update

Temporal Difference Learning

$$v_{t+1}(S_t) = v_t(S_t) + \alpha(R_{t+1} + \gamma v_t(S_{t+1}) - v_t(S_t))$$



Target

Temporal Difference Learning

$$v_{t+1}(S_t) = v_t(S_t) + \alpha(R_{t+1} + \gamma v_t(S_{t+1}) - v_t(S_t))$$



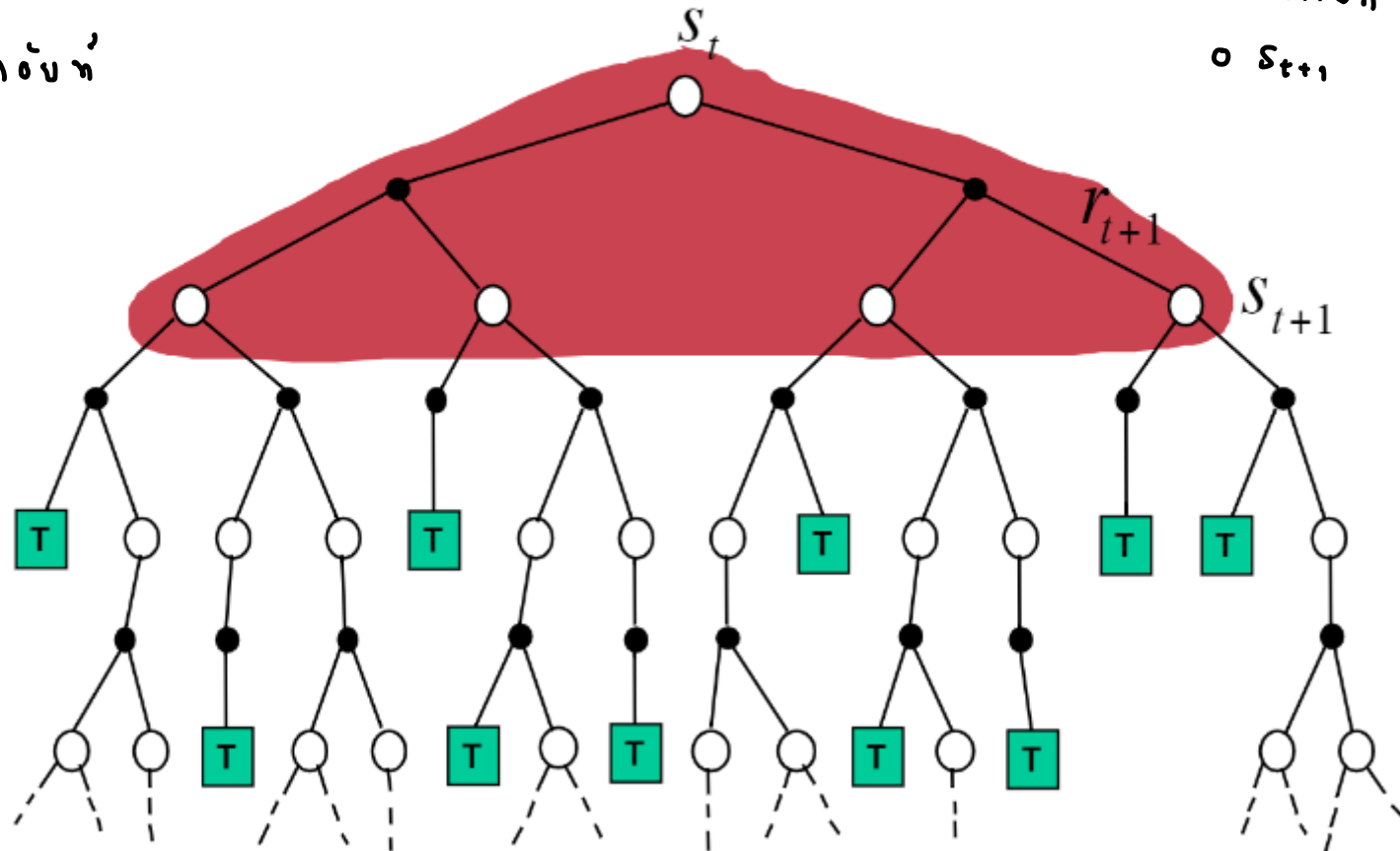
TD Error

Dynamic Programming Backup

$$v(S_t) \leftarrow \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) \mid A_t \sim \pi(S_t)]$$

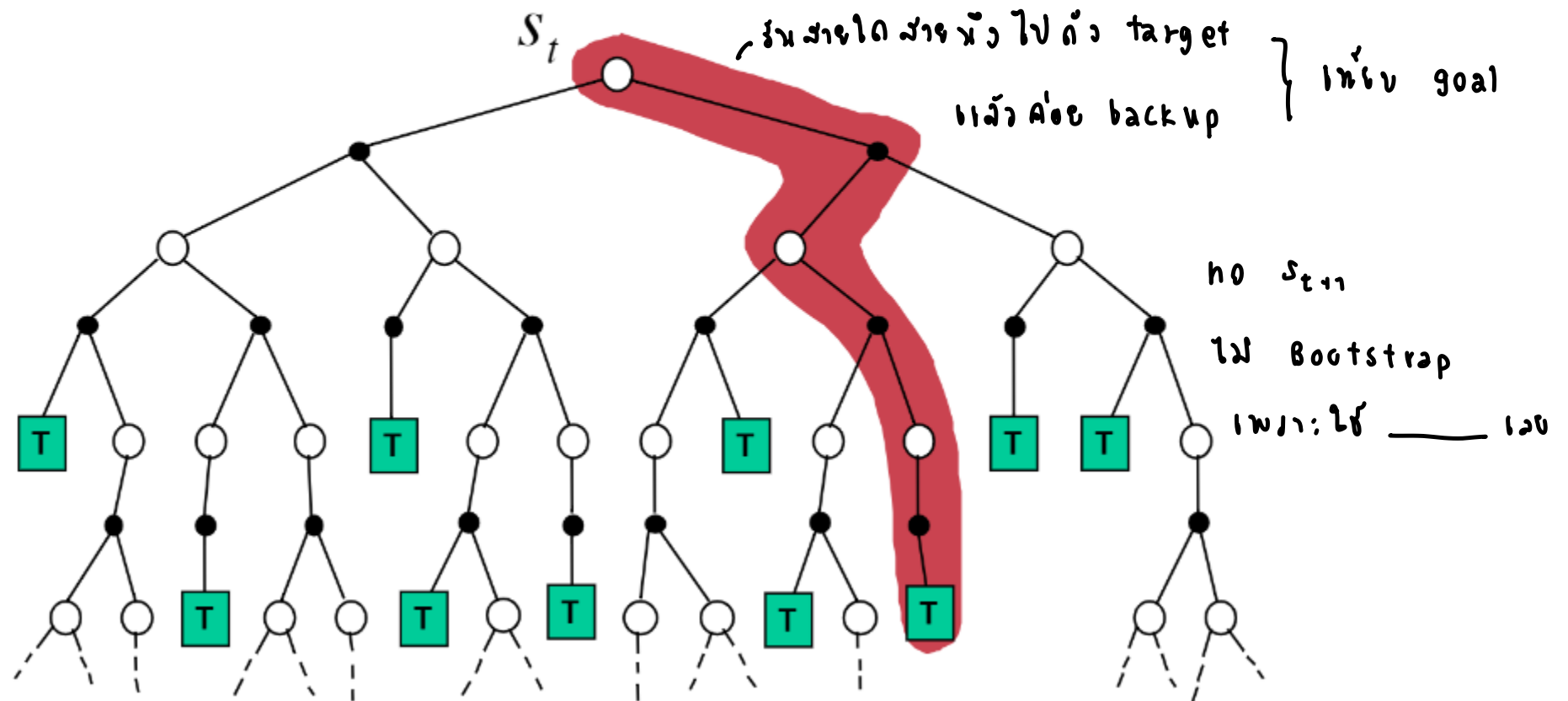
ทําคงส่วนที่นอก
ไว้ให้ *

- Action
- S_{t+1}



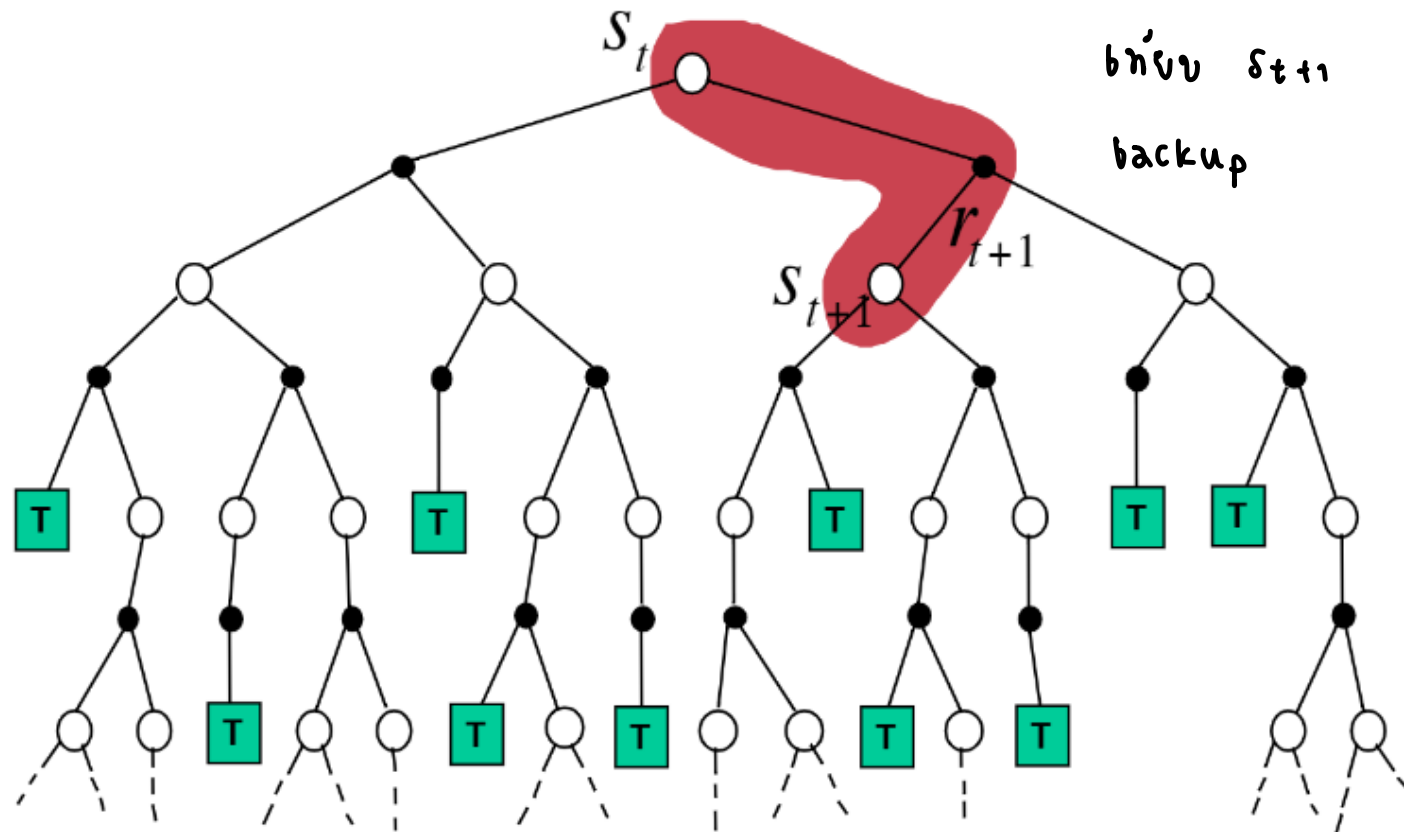
Monte-Carlo Backup

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



Temporal-Difference Backup

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



Bootstrapping and Sampling

Bootstrapping → การประมาณค่า

- Use the estimate of the next state to update
- DP and TD use
- MC does not

Sampling

- Update samples an expectation
- MC and TD sample
- DP does not

Temporal Difference Learning

We can also learn action values by updating value $q_t(S_t, A_t)$ with respect to estimated return $R_{t+1} + \gamma v(S_{t+1}, A_{t+1})$

$$q_{t+1}(S_t, A_t) = q_t(S_t, A_t) + \alpha_t(R_{t+1} + \gamma q_t(S_{t+1}, A_{t+1}) - q_t(S_t, A_t))$$

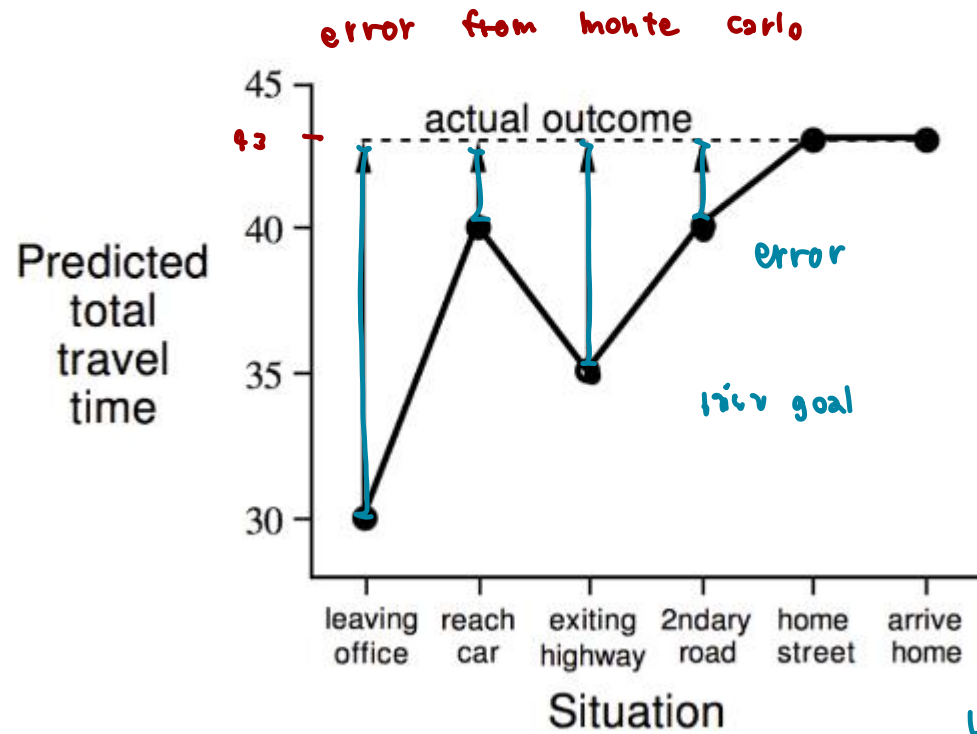
This algorithm is called **SARSA**, because it uses $S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}$

p_{TDL}

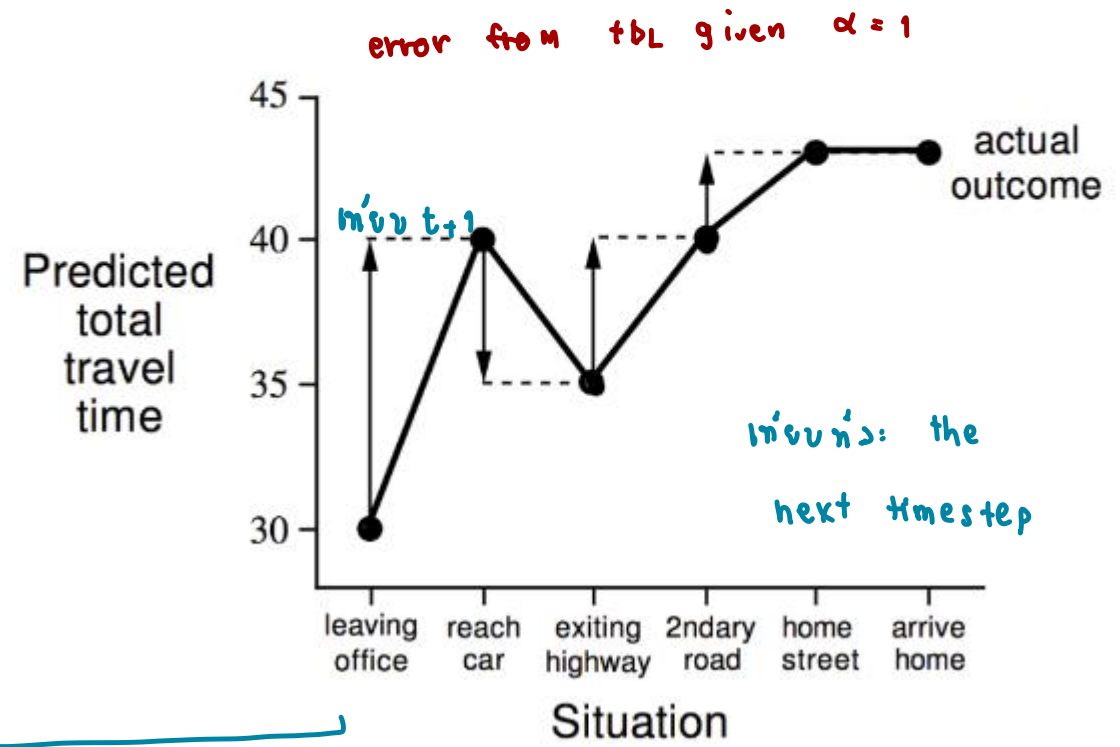
listen !!

State	Elapsed Time (minutes)	Predicted 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

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



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



TD vs MC

TD can learn **before** knowing the final outcome

- TD **learns** online **every step**.
- MC must wait until the end of the episode.

Continuous environment ដែល MC មិនអាចស្រាវជ្រាវបានទេ។

TD can learn **without** the final outcome ដោយ **terminate** មិនចាំបាច់

- TD can learn even if the full interaction sequence is incomplete. **MC must use the complete sequence**.
- TD can be used in a continuing environment. MC environment must terminate.

Bias-Variance Trade-off

→ bias free

Given MC return,

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

This is an unbiased estimation of $v_\pi(S_t)$

→ bias and no timestep delay

Given TD target,

$$R_{t+1} + \gamma v_t(S_{t+1})$$

This is a biased estimation of $v_\pi(S_t)$

MC : bias ↓ variance ↑

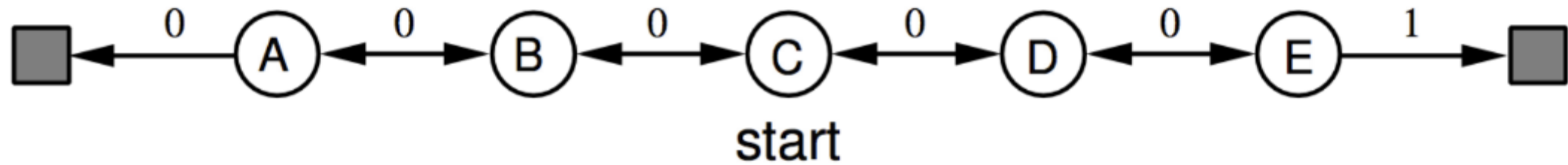
TD : " ↑ " ↓

Bias-Variance Trade-off

However, TD has **lower variance**

- MC return estimates from many random actions, transitions, and rewards
- TD target estimates from one random action, transition, and reward
- Variation of action, transition, and reward makes the TD target have less variance

Random Walk Example

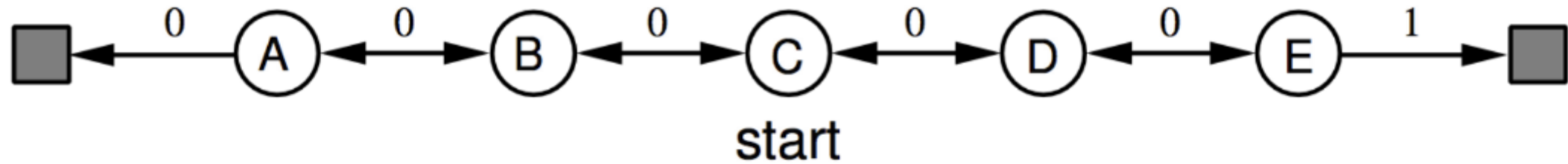


This game has 7 states, with uniform random transitions (50% left, 50% right)
 ทั่วๆไป

We can see that the true values of all states are:

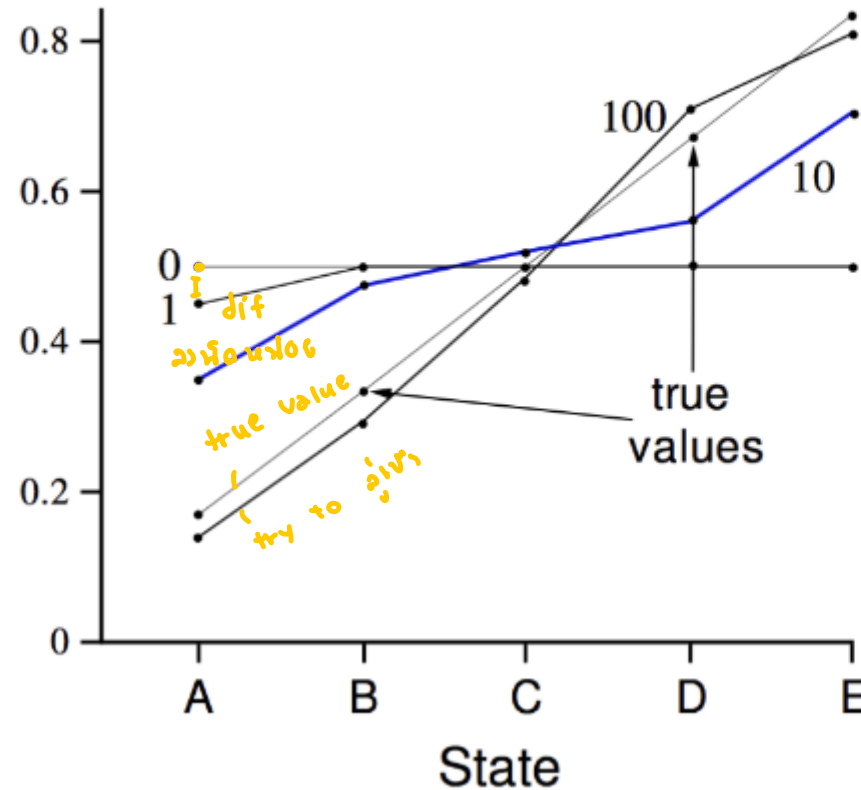
$$v(A) = \frac{1}{6}, v(B) = \frac{2}{6}, v(C) = \frac{3}{6}, v(D) = \frac{4}{6}, v(E) = \frac{5}{6}$$

Random Walk Example



Given that initial values:
 $v(s) = 0.5, \forall s$

Estimated
value

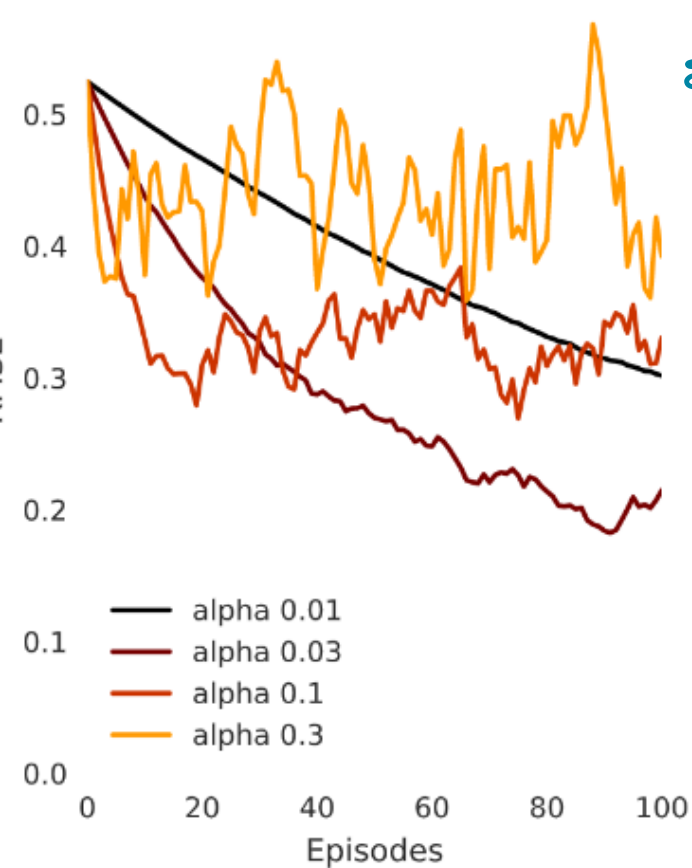
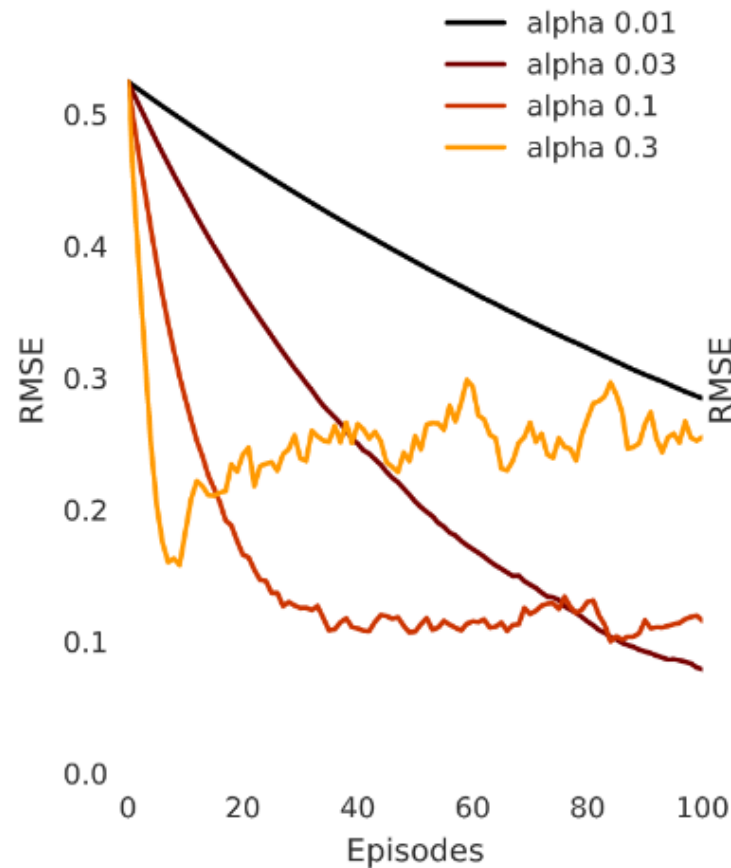


Random Walk Example

TD

MC

HW 2 :
try this



alpha ↑
oscillations

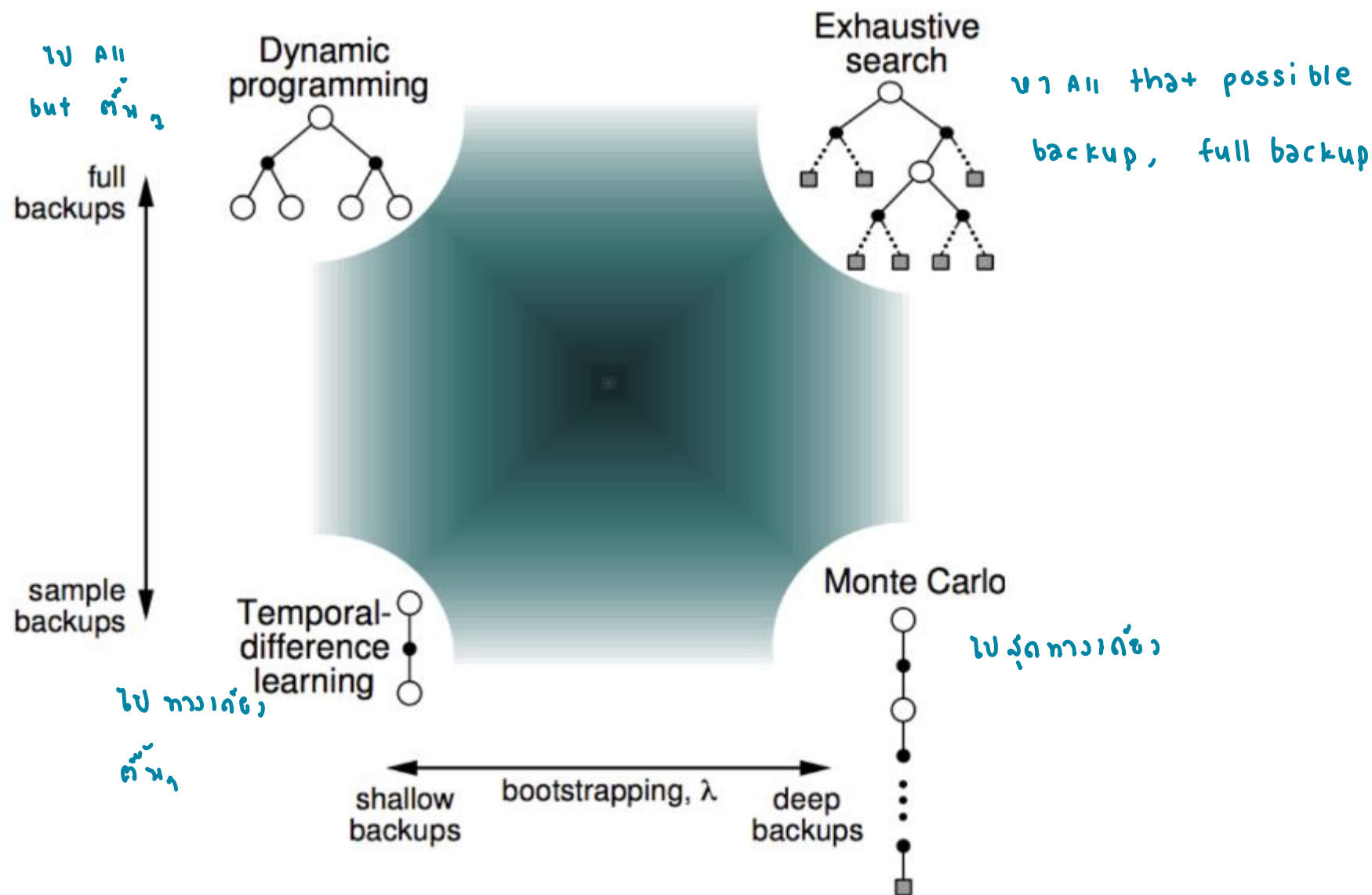
TD vs MC

TD exploits Markov property *each step*

- TD can help in fully-observable environments

MC does not exploit Markov Property *all step*

- MC can help in partially-observable environments



Multi-Step Updates

Multi-Step Updates

TD and MC tackle the problem from different angles

- TD estimates, and bootstrap
- MC use true return, but is noisy

Can we do something between?

- Multi-Step Prediction ★

Multi-Step Returns

We can arrange TD and MC as:

TD $n = 1$ $G_t^{(1)} = R_{t+1} + \gamma v(S_{t+1})$ ↖ อยากรู้ว่า: กับ 1 : want to predict 1 step ข้างหน้า

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

\vdots ↖ sample 2 timestep แล้วค่อย Bootstrap

MC $n = \infty$ $G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T$ ↖ sample หมด (feel like MC)

In general, we can formalize n-step return as

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

Multi-Step Returns

In general, we can formalize n-step return as

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

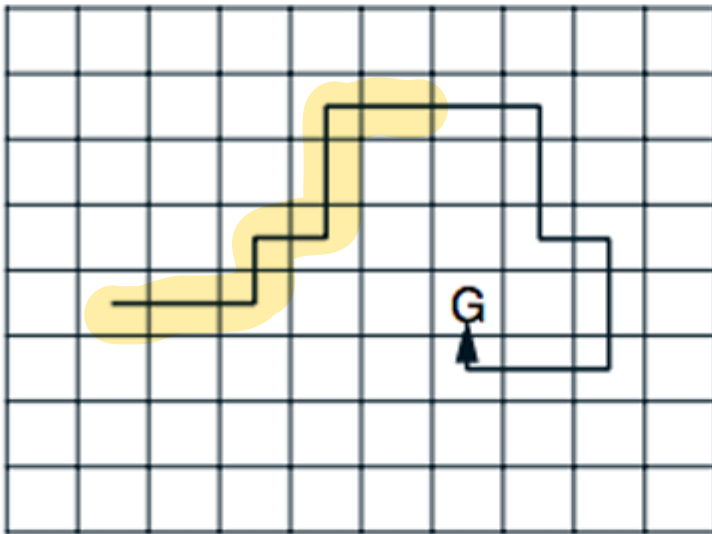
We can define multi-step temporal-difference learning as

$$v(S_t) \leftarrow v(S_t) + \underbrace{\alpha (G_t^{(n)} - v(S_t))}_b$$

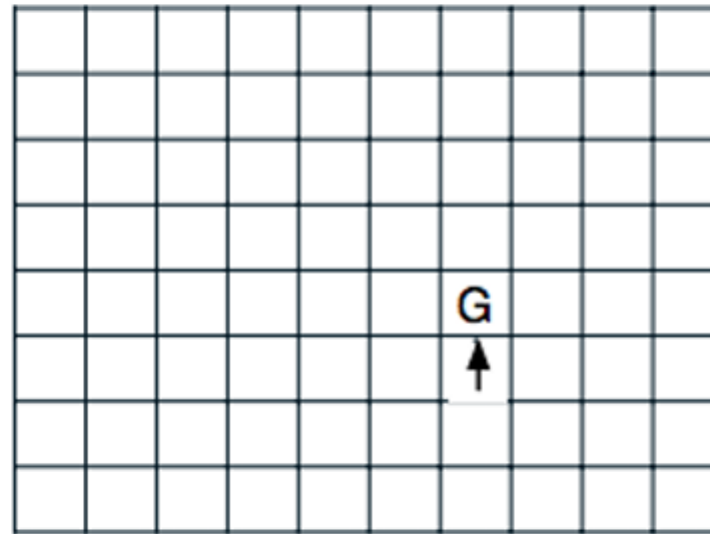
if $G_t = v(S_t)$ and $\alpha > 0$ then $\Delta v(S_t) = \alpha (G_t - v(S_t))$
 b
 TD's target - current value

Example

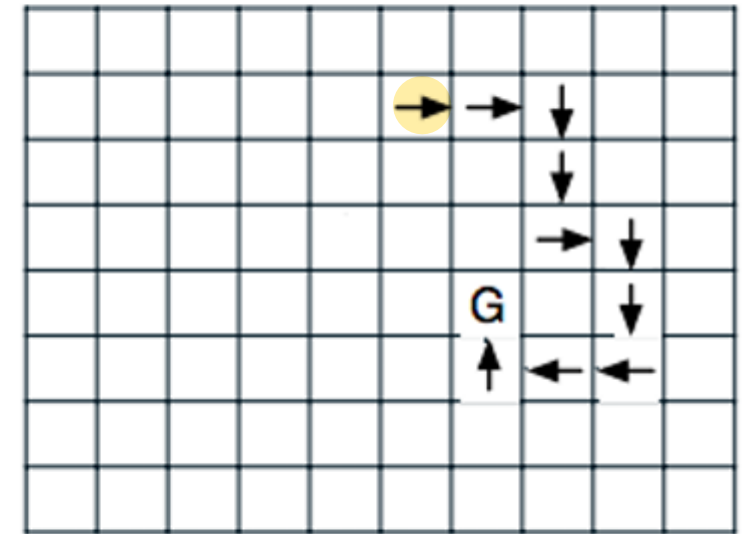
Path taken



Action values increased by one-step Sarsa

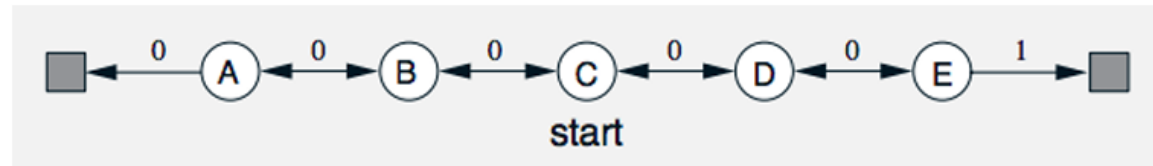


Action values increased by 10-step Sarsa

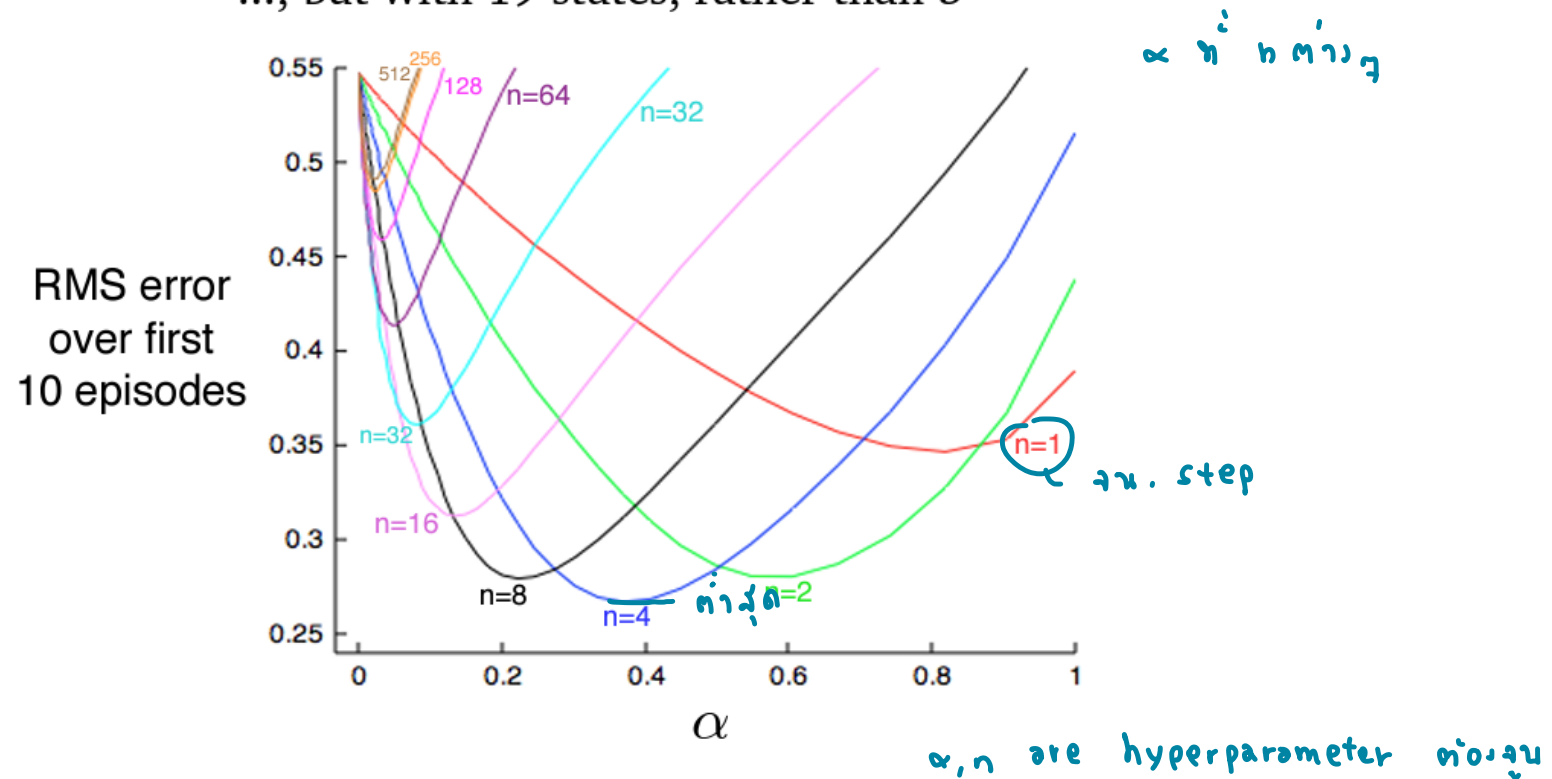


(Reminder: SARSA is TD for action values $q(s, a)$)

Example



..., but with 19 states, rather than 5



Mixing Multi-Step Returns

Consider a multi-step return equation,

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

Multi-step returns bootstrap on one state, $v(S_{t+n})$. However, we can modify the equation so you can bootstrap a little bit on many states

$$G_t^\lambda = R_{t+1} + \gamma((1-\lambda)v(S_{t+1}) + \lambda G_{t+1}^\lambda)$$

reward + decaying factor $(1-\lambda)v(S_{t+1})$ + decay (return)

๒ จงวนการ

Multi-Step Returns

$$G_t^\lambda = R_{t+1} + \gamma((1 - \lambda)v(S_{t+1}) + \lambda G_{t+1}^\lambda)$$

hyperparameter λ

If we consider:

$$\lambda = 0 \quad G_t^{\lambda=0} = R_{t+1} + \gamma v(S_{t+1}) \quad (\text{TD})$$

$$\lambda = 1 \quad G_t^{\lambda=1} = R_{t+1} + \gamma G_{t+1} \quad (\text{MC})$$

Multi-Step Returns

balance

Using multi-step returns,

- We can utilize benefits of TD and MC
- Bootstrapping may cause bias
- Monte Carlo involves high variance

Generally, using multi-step returns are good

Summary

No dynamic model of problem
cannot find expectation value
then we sampling by using
Monte Carlo run all to learn
and TD run step+1 to learn
but they are all have curse
then we merge MC & TD to
multi-step returns.