

Reinforcement Learning

Lecture 5

Model-free Prediction

Monte Carlo

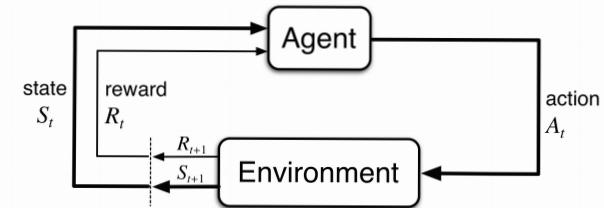
Temporal Differences

Robert Peharz

Institute of Theoretical Computer Science
Graz University of Technology

Winter Term 2023/24

Recap: Basic Notions in RL



Markov Decision Process (MDP) (S, A, P)

state space S , action space A , dynamics $p(s', r | s, a)$

Policy $\pi(a|s)$

Discounted Return

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}, \quad \gamma \in [0, 1]$$

Value Function

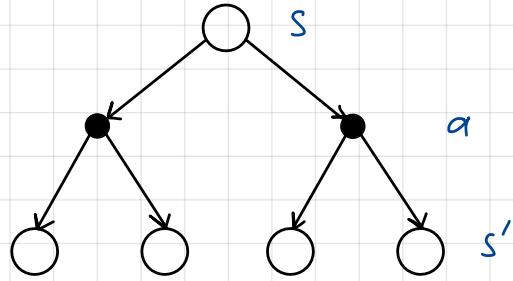
$$V_\pi(s) := E_\pi[G_t \mid S_t = s] = E_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

q-Function (state-action value function)

$$q_\pi(s, a) := E_\pi[G_t \mid S_t = s, A_t = a] = E_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

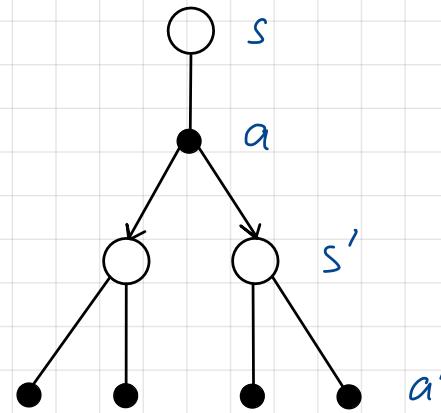
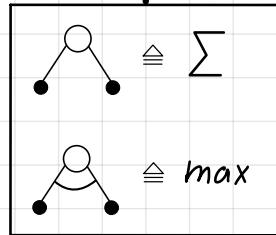
Recap: Bellman Equations, Backup Diagrams

Bellman Expectation Equation in v



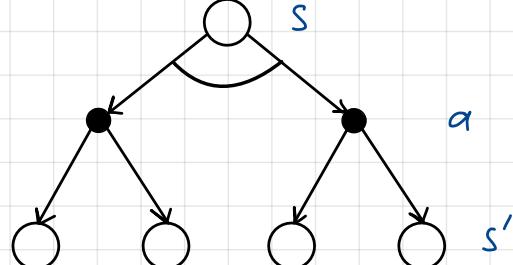
$$v_{\pi}(s) = r(s) + \gamma \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) v_{\pi}(s')$$

Bellman Expectation Equation in q



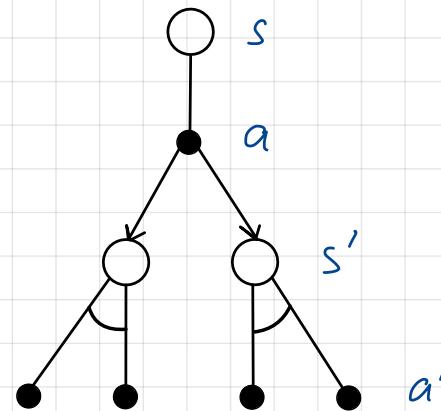
$$q_{\pi}(s,a) = r(s,a) + \gamma \sum_{s'} p(s'|s,a) \sum_{a'} \pi(a'|s') q_{\pi}(s',a')$$

Bellman Optimality Equation in v



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

Bellman Optimality Equation in q



$$q_{\pi}(s,a) = r(s,a) + \gamma \sum_{s'} p(s'|s,a) \max_{a'} q_{\pi}(s',a')$$

Recap: Policy Improvement, Policy Iteration, Value Iteration

greedy policy given current policy π

$$\pi'(a|s) = \begin{cases} 1 & \text{if } a = \arg \max_{a'} q_\pi(s, a') \\ 0 & \text{otherwise} \end{cases}$$

$$V^*(s) := \max_{\pi} V_{\pi}(s)$$

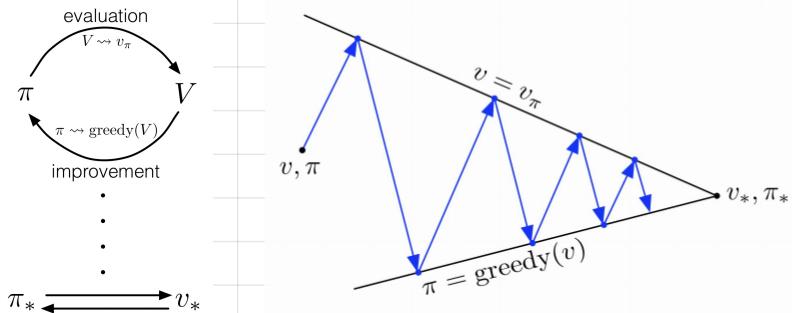
$$q^*(s, a) := \max_{\pi} q_{\pi}(s, a)$$

Policy Improvement Theorem

$$V_{\pi}(s) \leq V_{\pi'}(s) \quad \forall s \in S$$

Theorem: Global Optimal Policy

For any initial policy π_0 , policy iteration converges to an optimal policy with value function v^*



Value Iteration

Policy Iteration with Truncated Policy Evaluation

Policy Iteration

1 number k of policy evaluation steps $\dots \infty$

Model-free Prediction, Monte Carlo

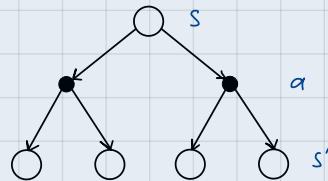
So far: Known MDP Dynamics

Planning

Bellman

Expectation Equation

$$p(s', r | s, a)$$

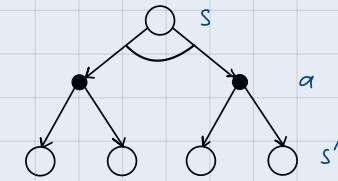


$$V_\pi(s) = r(s) + \gamma \mathbb{E}[V_\pi(s_{t+1})]$$

Policy Evaluation

Bellman

Optimality Equation



$$V_\pi(s) = \max_a r(s, a) + \gamma \mathbb{E}[V_\pi(s_{t+1}) | a]$$

(Generalized) Policy Iteration

Today: Unknown MDP Dynamics

Reinforcement Learning

We can, however, interact with MDP and produce samples (experience)

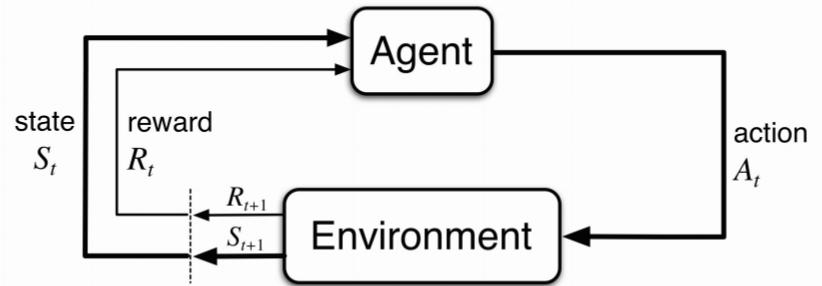
Model-free RL directly learn $\tilde{\pi}$ ($v_{\tilde{\pi}}$, $q_{\tilde{\pi}}$)

Today

Model-based RL learn (some aspect of) MDP as well

We start with prediction / policy evaluation (compute $v_{\tilde{\pi}}/q_{\tilde{\pi}}$, $\tilde{\pi}$ fixed)

Model-free Policy Evaluation



- MDP dynamics $p(s', r | s, a)$ (environment) unknown
- Policy $\pi(a | s)$ given
- We can collect experience in the form of episodes

$$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, \dots, S_K \sim \hat{\pi}, p$$

i.e., assume episodic tasks for now

How to compute the value function?

$$V_{\hat{\pi}}(s) = \mathbb{E}_{\hat{\pi}}[G_t \mid S_t = s] = \mathbb{E}_{\hat{\pi}} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

Monte Carlo (MC) Estimation

state space of \underline{X}

Let \underline{X} be a random vector and $f: \underline{\mathcal{X}} \rightarrow \mathbb{R}$ be a function.
Assume that we want to compute $\mathbb{E}[f(\underline{X})]$ and that we have
N i.i.d. (independent and identically distributed) samples, i.e.

$$\underline{X}_i \sim p_{\underline{X}}, \quad 1 \leq i \leq N$$

The Monte Carlo estimator for $\mathbb{E}[f(\underline{X})]$ is defined as

$$\hat{E} := \frac{1}{N} \sum_{i=1}^N f(\underline{X}_i)$$

(Note that \hat{E} is a random variable)

\hat{E} is unbiased, i.e.

$$\mathbb{E}[\hat{E}] = \mathbb{E}[f(\underline{X})]$$

Its variance is

$$\text{var}(\hat{E}) = \frac{\text{var}(f(\underline{X}))}{N}$$

Thus, $\hat{E} \xrightarrow{N \rightarrow \infty} \mathbb{E}[f(\underline{X})]$

Estimating Area of Unit Circle with MC

- What is the area of the unit circle? Easy, $\pi^2 \approx \pi = 3.1415\ldots$
- One way to check this empirically is Monte Carlo:
 - define uniform distribution on unit square $[-1, 1] \times [-1, 1]$

$$p(x, y) = \begin{cases} 1/4 & \text{if } |x| \leq 1 \text{ and } |y| \leq 1 \\ 0 & \text{o.w.} \end{cases}$$

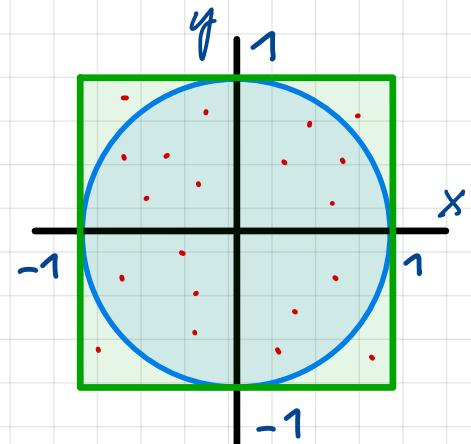
- area of unit circle:

indicator function $\mathbb{I}[a] = \begin{cases} 1 & \text{if } a \text{ is true} \\ 0 & \text{o.w.} \end{cases}$

$$\begin{aligned} \iint_{-1}^1 \mathbb{I}[x^2 + y^2 \leq 1] dx dy &= \iint 4 p(x, y) \mathbb{I}[x^2 + y^2 \leq 1] dx dy \\ &= 4 \mathbb{E}[\mathbb{I}[x^2 + y^2 \leq 1]] \end{aligned}$$

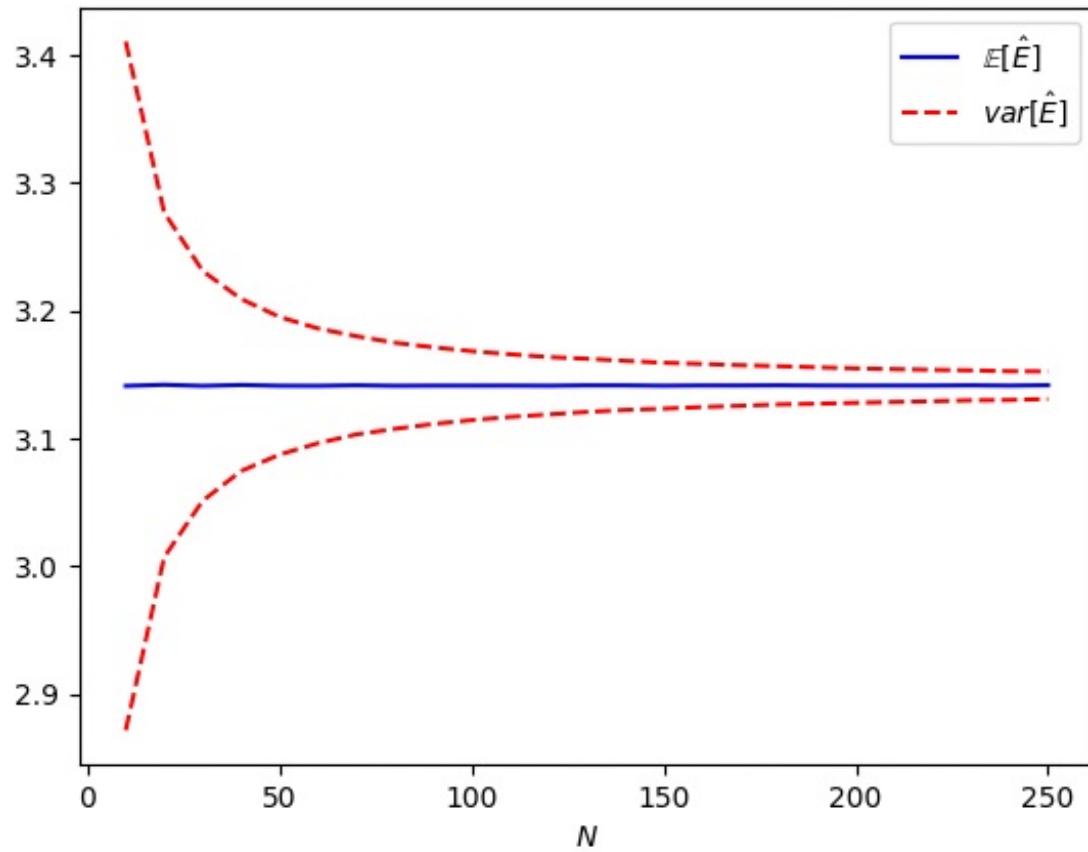
- simulate N draws (X_i, Y_i) from unit square

$$\hat{\pi} \approx 4 \cdot \hat{E} = 4 \cdot \frac{1}{N} \sum_{i=1}^N \mathbb{I}[X_i^2 + Y_i^2 \leq 1]$$



Estimating Area of Unit Circle MC con't

Mean and variance of Monte Carlo estimator $4\hat{E}$



(estimated by 10^6 repeated Monte Carlo runs, for each N)

Monte Carlo

- ✓ unbiased
- ✓ asymptotically exact
- ✓ P_{π} not required
- ✓ { if not (explicitly) required } just needs to be simulated

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

- simulate many episodes
- compute G
- average
- efficiency: also use sub-episodes for respective starting states

$S_0, A_0, R_1, S_1, A_1, R_2, \overbrace{S_2, A_2, R_3, \dots, S_{T-1}, A_{T-1}, R_T}^{\text{---}}$

Monte Carlo Evaluation (Prediction)

estimate with Monte Carlo

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t | S_t = s] = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right]$$

MC Prediction

- Given: π
- Initialize empty list $\text{Returns}(s)$ for all $s \in \mathcal{S}$
- Initialize $v(s) = 0$ for all $s \in \mathcal{S}$

repeat

- generate an episode $S_0, A_0, R_1, \dots, S_{T-2}, A_{T-2}, R_{T-1}, S_{T-1}, A_{T-1}, R_T$

$$G \leftarrow 0$$

- iterate $t = T-1, \dots, 0$:

$$G \leftarrow \gamma G + R_{t+1}$$

- append G to $\text{Returns}(S_t)$

$$v(S_t) = \text{average}(\text{Returns}(S_t))$$

// this recursively computes
 $G_t = \sum_k \gamma^k R_{t+k+1}$ for each sub-episode

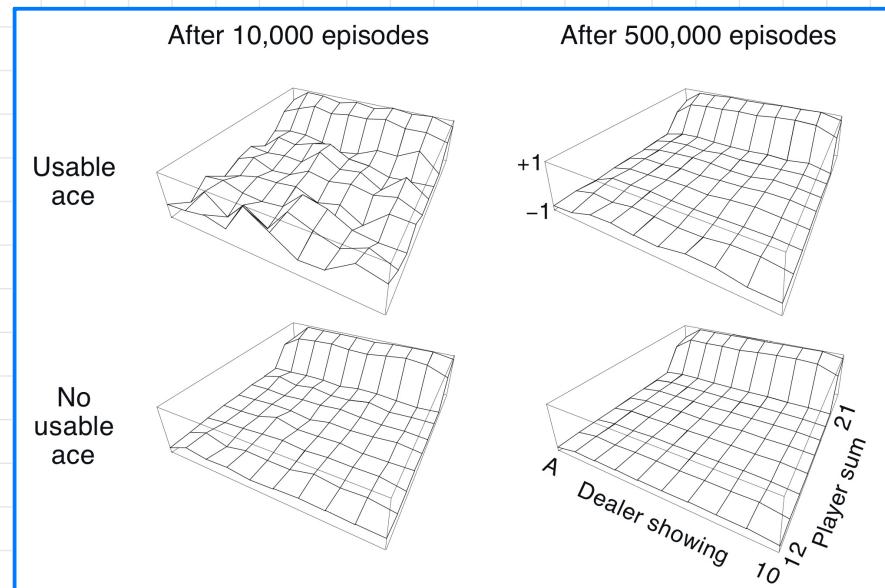
Black Jack

(Sutton & Barto, p. 93)

- you and the dealer get 2 cards
- one of the dealer's cards is face up (visible)
- you have the following options (actions)
 - hit: ask for another card
 - stick: stop and use your current hand
- if you get more than 21 points you are **bust** (you lose)
- dealer plays after you: they stick with 17 or higher
- win: $R=1$, lose: $R=-1$, draw: $R=0$
- assume policy: stick with 20 or 21, otherwise hit

2:	2 points
3:	3 points
:	
10:	10 points
face cards:	10 points
ace:	1 or 11 points

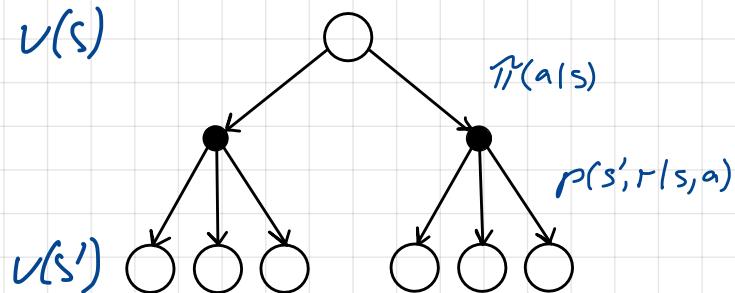
- easy to simulate
- Monte Carlo results →



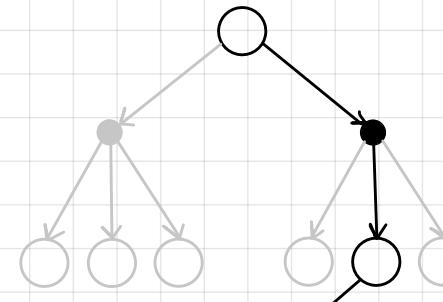
Temporal Difference Prediction

Unifying Planning and Monte Carlo?

Planning (Dynamic Programming)



Monte Carlo



What's in between?

- one step look-ahead
considers all options

$$\text{IE}[\dots] = \sum_{a,s,r} \dots$$

- uses current estimate
 $v(s_{t+1})$ ($q(s_{t+1}, a_{t+1})$)
as "backup" (Bellman equations)

- long range look-ahead
- considers one option for each t
- uses empirical G_t samples

Intermezzo: Online Monte Carlo

- Monte Carlo

- maintain a set $\text{Returns}(s)$ for each $s \in S$

- $V_{\pi}(s) = E_{\pi}[G_t | s] \approx \text{average}(\text{Returns}(s))$

- Online Monte Carlo

- $v(s) \leftarrow 0, N(s) \leftarrow 0$

- whenever a new sample G_t becomes available for s :

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

$$v(s) \leftarrow v(s) + \frac{1}{N(s)} (G_t - v(s))$$

- generalize this to a learning rule

$$v(s) \leftarrow v(s) + \alpha (\overbrace{G_t - v(s)}^{\text{"\$"}})$$

step size target current estimate

- constant α : online learning, forgetting long past data

online (incremental) averages:

$$\begin{aligned}\bar{x}_N &= \frac{1}{N} \sum_{i=1}^N x_i \\ &= \frac{1}{N} \left(x_N + \sum_{i=1}^{N-1} x_i \right) \\ &= \frac{1}{N} (x_N + (N-1) \bar{x}_{N-1}) \\ &= \frac{1}{N} (x_N + N \bar{x}_{N-1} - \bar{x}_{N-1}) \\ &= \bar{x}_{N-1} + \frac{1}{N} (x_N - \bar{x}_{N-1})\end{aligned}$$

Temporal Difference (TD) Prediction

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

target

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$= R_{t+1} + \gamma G_{t+1}$$



Thus, could $R_{t+1} + \gamma V(s')$ be a good target? Still correct algorithm?

Temporal Difference (TD) Prediction

- Initialize $v(s)$ arbitrarily, except $v(s) = 0$ for terminal s
 - for each episode
 - init s
 - repeat
 - take action $a \sim \pi(a|s)$, observe s', r
 - $v(s) \leftarrow v(s) + \alpha (r + \gamma v(s') - v(s))$
 - $s \leftarrow s'$

// for continuing task,
this is one "infinite episode"

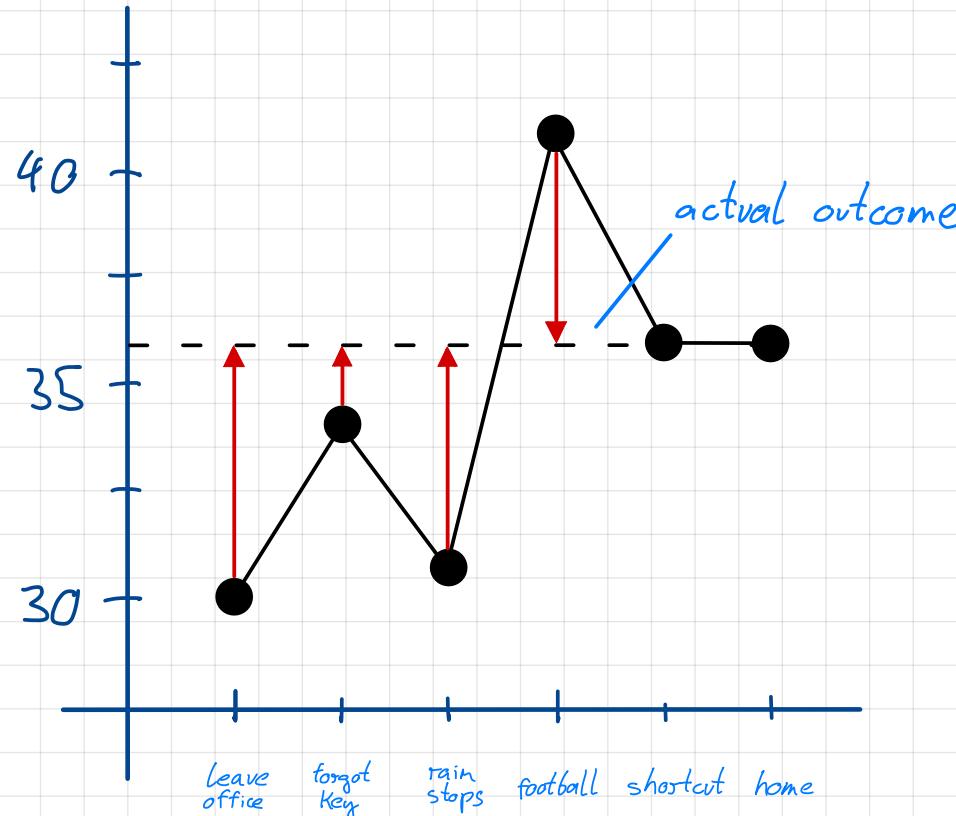
Cycling Home (from a Dutch University)

state	elapsed time	time updates	predicted total time
leaving office	0	0	30
forgot bike key on desk	2	+4	34
Leave campus, rain starts	6	-3	31
local football club has won against Ajax Amsterdam: drunk fans block bike path	15	+10	41
construction of shortcut bike path completed	30	-5	36
arriving home	36	0	36

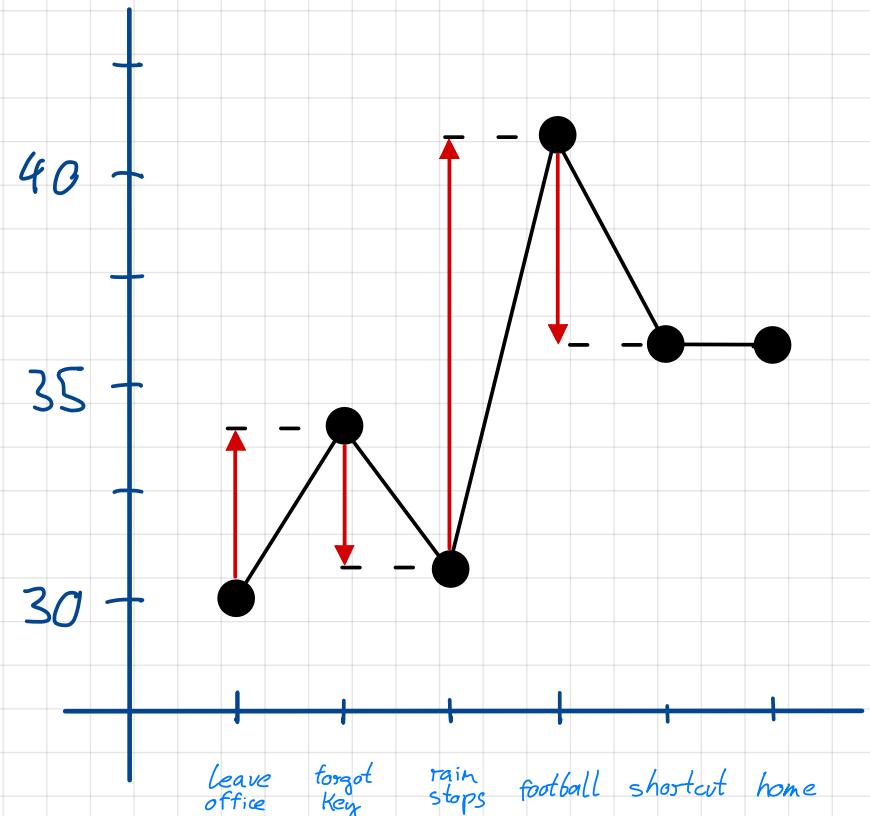


adapted from Sutton & Barto, p. 122

Cycling Home (from a Dutch University)



Updates via Monte Carlo



Updates via Temporal Differences

Soundness of TD

$$v(s) \leftarrow v(s) + \alpha \underbrace{(R_{t+1} + \gamma v(s') - v(s))}_{\delta_t : \text{TD error}}$$

Assume $\mathbb{E}[\delta_t | S_t = s] = 0$

Thus, $\mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) - v(S_t) | S_t = s] = 0$

$$\mathbb{E}[v(S_t) | S_t = s] = \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) | S_t = s]$$

$$v(s) = \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) | S_t = s] \quad \text{Bellman expectation equation} !$$

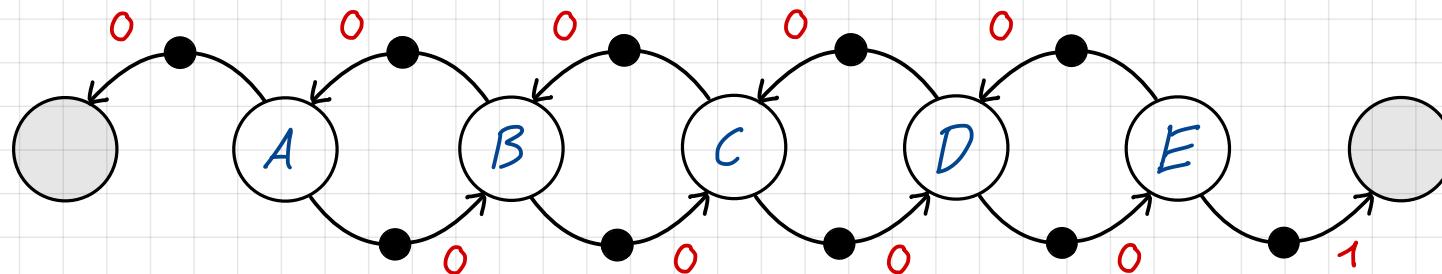
Intuitively, TD aims to establish Bellman expectation equation, at which point $v = v_{\pi}$. This indeed happens in the infinite data limit if stepsize α_t is set correctly:

$$\sum_{t=0}^{\infty} \alpha_t = \infty, \sum_{t=0}^{\infty} \alpha_t^2 < \infty \quad \text{e.g. } \alpha_t = \frac{1}{t}$$

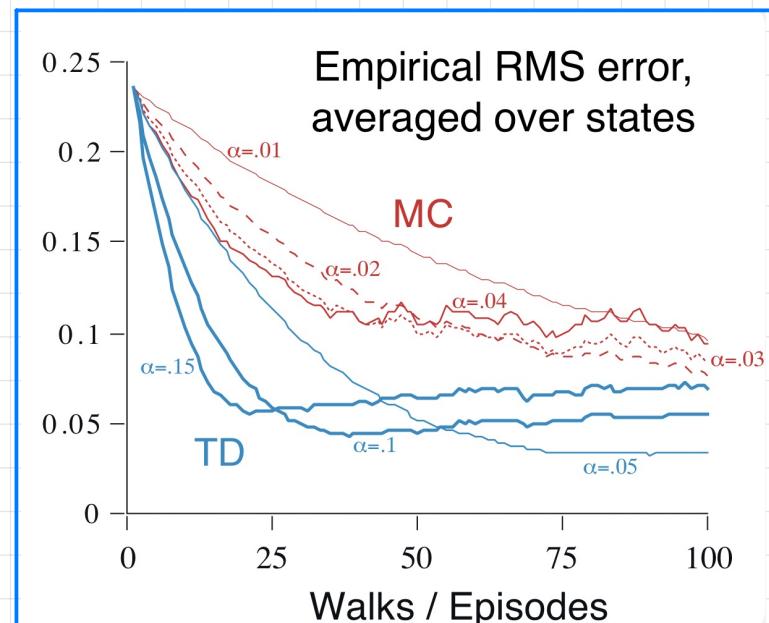
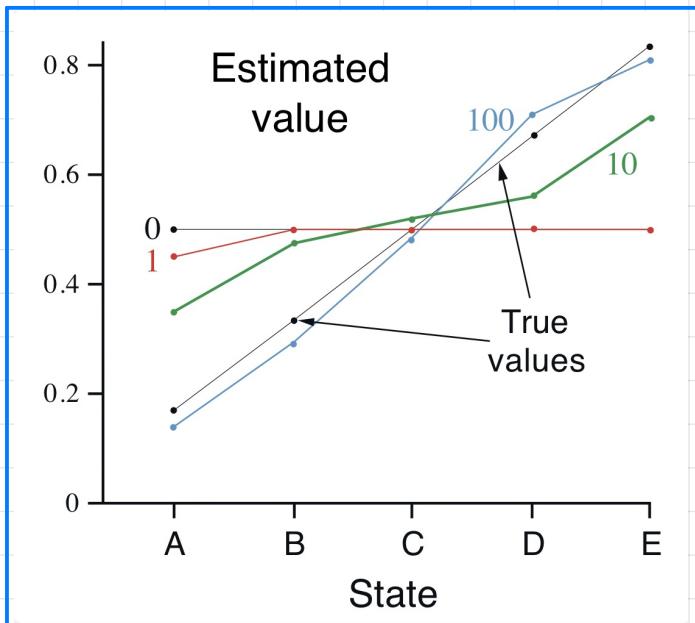
(standard conditions from stochastic optimization, cf. Robins-Monroe, 1951)

Monte Carlo vs. TD

(Sutton & Barto, p.125)



$$\hat{\pi}(a|s) = 0.5$$



Value function learned by TD

MC vs. TD

Monte Carlo vs. TD II

(Sutton & Barto, p. 127)

How do MC and TD compare on finite data?

Consider the following data for state and reward (ignoring actions)

1: A, 0, B, 0

2: B, 1

3: B, 1

4: B, 0

5: B, 1

6: B, 1

7: B, 1

8: B, 1

How would you estimate $v(A)$ and $v(B)$?

Assume undiscounted return.

Monte Carlo vs. TD II

(Sutton & Barto, p. 127)

For B, it is sensible to assume $v(B) = \frac{6}{8} = \frac{3}{4}$, since in 6 out of 8 cases when being in state B we got 1 reward, 0 otherwise.

There are two meaningful answers for A

1) the "Monte Carlo way"

once observed A and got $G=0$, thus, estimate $v(A)=0$.

2) the "TD way"

whenever we were in A, we got 0 reward and went to B, which we estimated as $v(B) = \frac{3}{4}$. Thus, estimate $v(A) = \frac{3}{4}$.

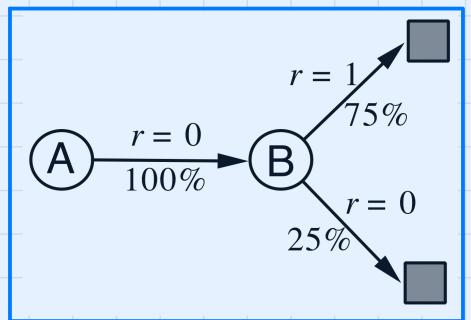
TD is equivalent to

1) fitting the MDP with maximum likelihood

2) solving MDP with Dynamic Programming

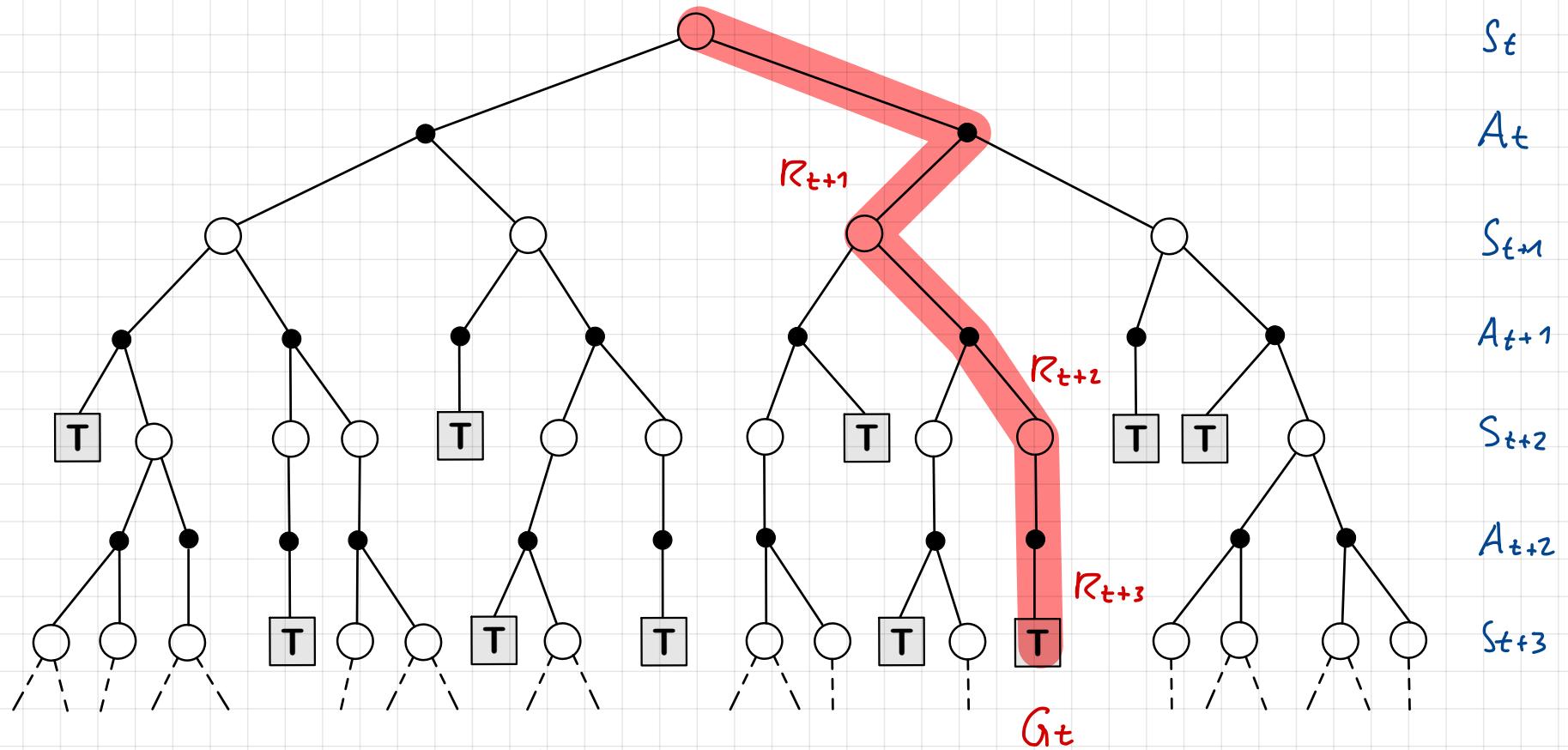
MC does not exploit the MDP structure.

In fact, MC still works when Markov property is violated!



Monte Carlo Backup

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

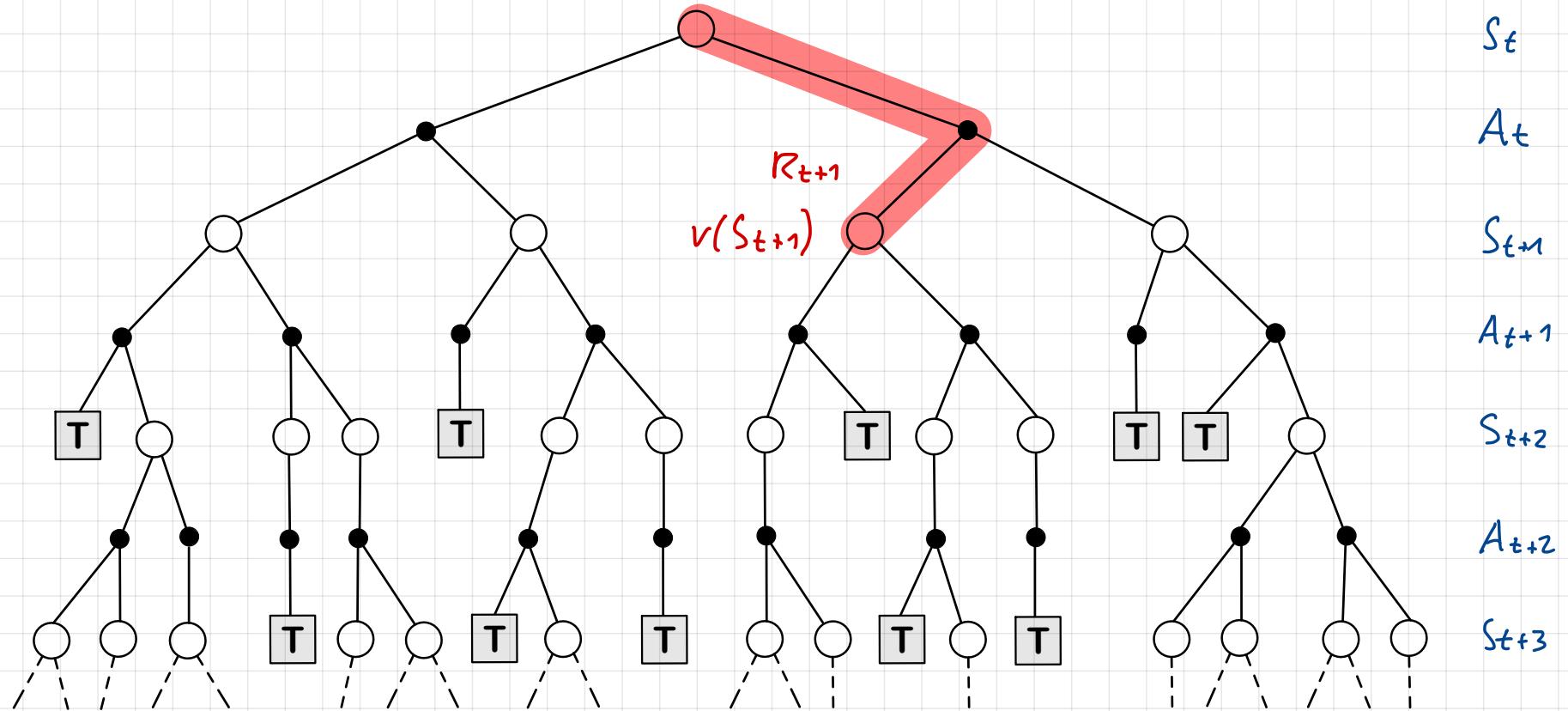


Adapted from D. Silver

TD Backup

uses current estimate v
as part of target: "bootstrapping" *

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

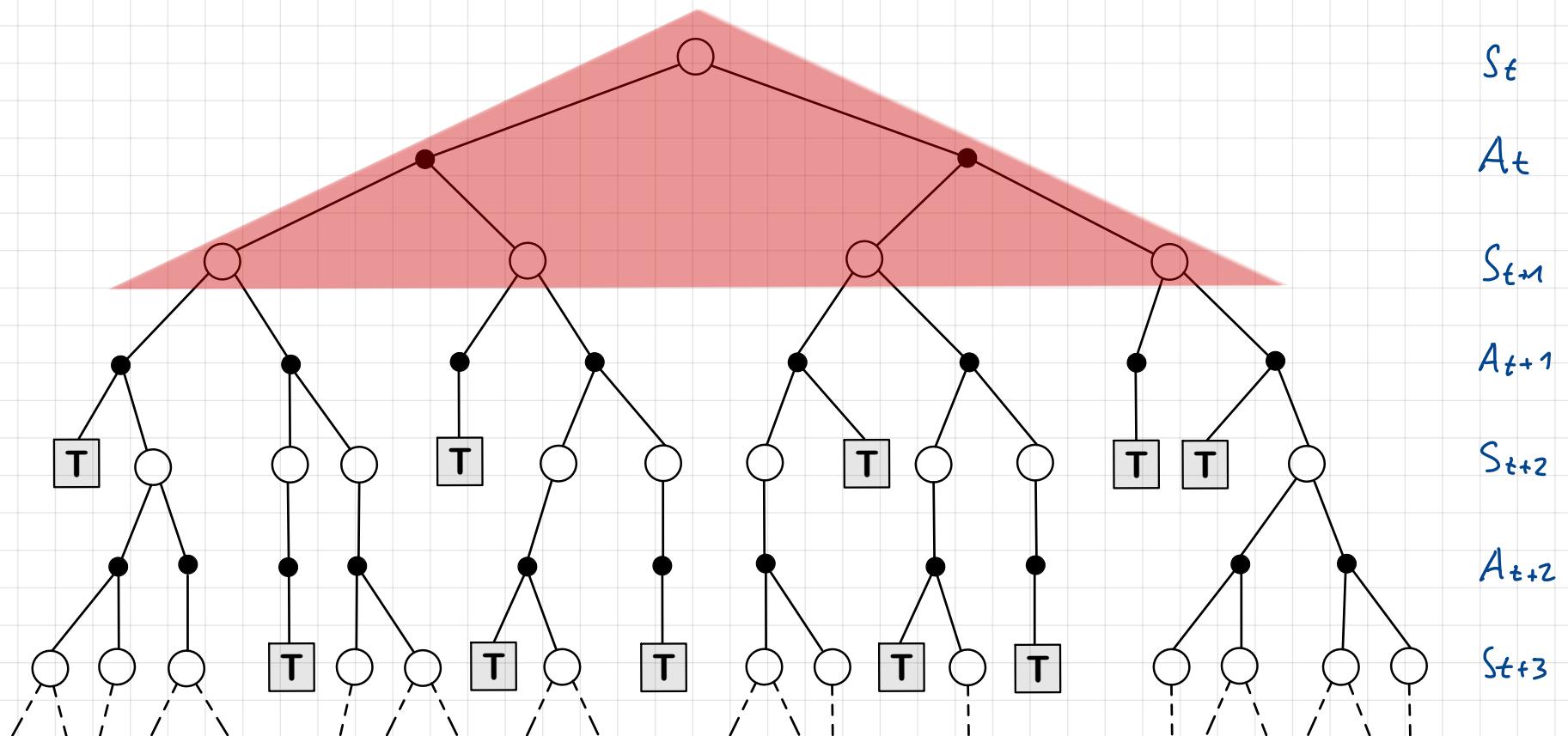


* don't confuse with bootstrap in statistics

Dynamical Programming Backup

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

bootstrapping



Properties of Prediction Algorithms

	allows for continuing tasks	bootstraps (uses current v as proxy)	model-free	samples	Look-ahead
Dynamic Programming	✓	✓			1
Monte Carlo			✓	✓	∞
Temporal Differences	✓	✓	✓	✓	1

N-step Time Differences

We can easily come up with schemes in between TD and MC.
General sampling-based update rule

$$v(s) \leftarrow v(s) + \alpha (\hat{g} - v(s))$$

target

$$(1\text{-step}) \text{ TD: } \hat{g} = R_{t+1} + \gamma v(S_{t+1})$$

$$2\text{-step TD: } \hat{g} = R_{t+1} + \gamma R_{t+2} + \gamma^2 v(S_{t+2})$$

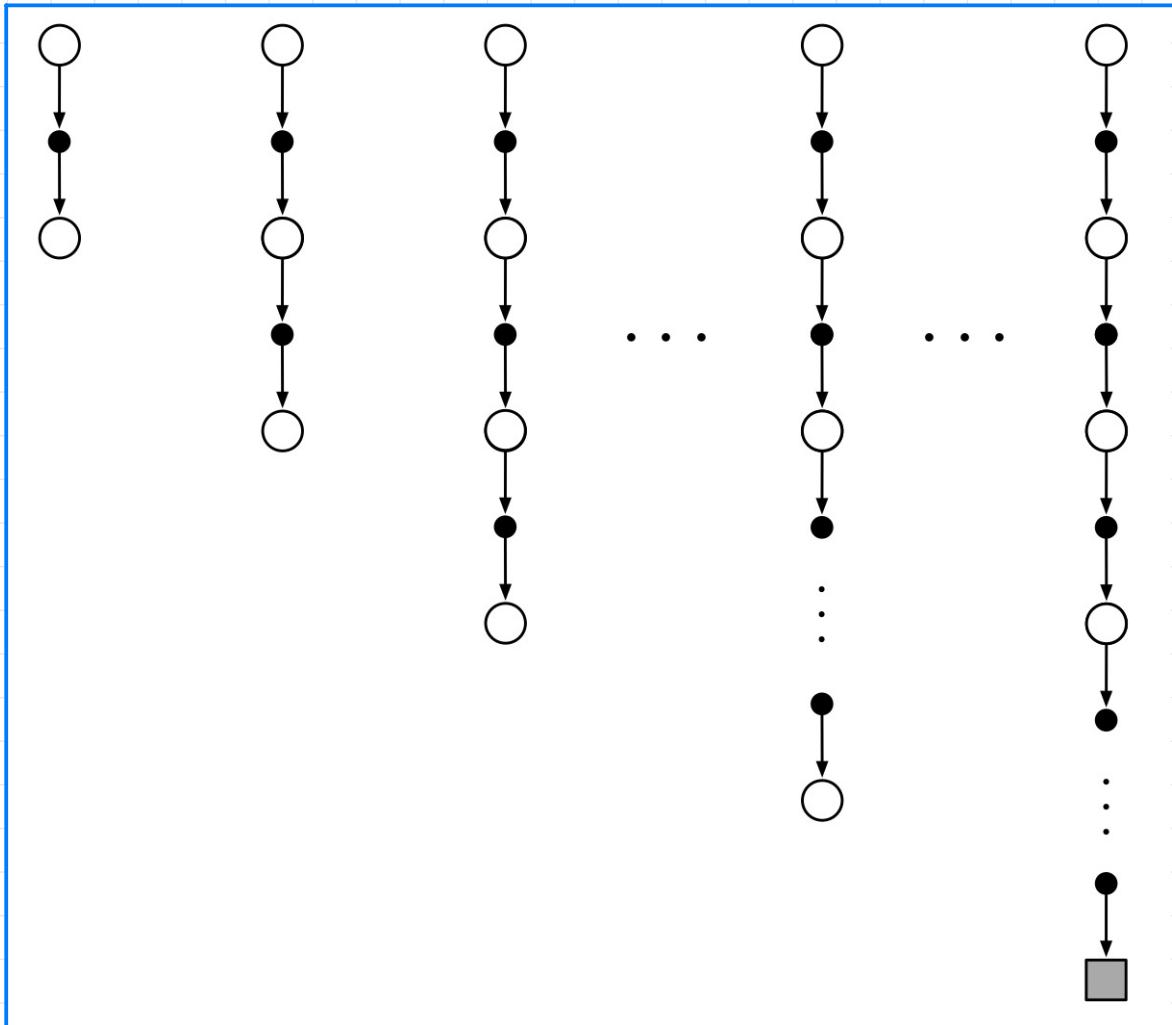
$$3\text{-step TD: } \hat{g} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 v(S_{t+3})$$

⋮

$$\infty\text{-step TD (MC)} \quad \hat{g} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots = g_t$$

N-step Time Differences

1-step TD 2-step TD 3-step TD n-step TD ∞ -step TD
MC



Unified View of RL Algorithms

