

# Lecture 4: Monte Carlo Methods and Temporal-Difference Learning

Zhi Wang & Chunlin Chen

Department of Control and Systems Engineering  
Nanjing University

Oct. 25th, 2021

# Table of Contents

## 1 Monte Carlo Methods

- Monte Carlo prediction
- Monte Carlo control

## 2 Temporal-Difference Learning

- Temporal-difference prediction
- Temporal-difference control

## 3 Extensions

- $n$ -step bootstrapping
- Eligibility traces

# Model-based vs. Model-free

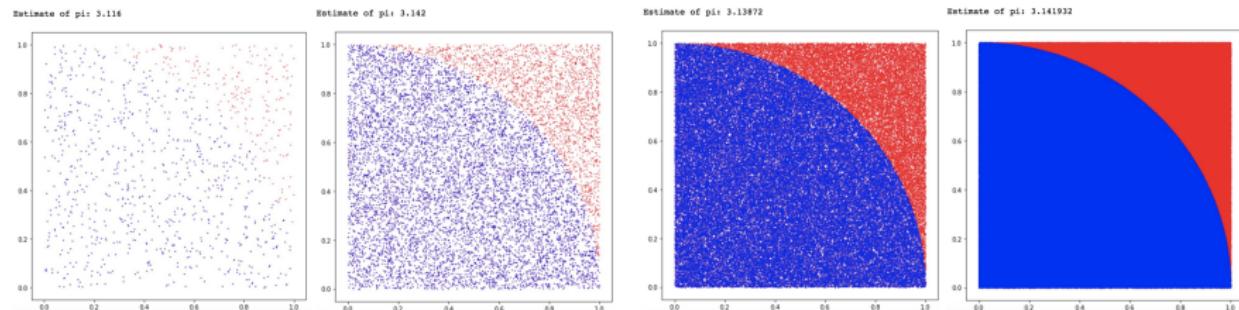
- Model-based algorithms, e.g., dynamic programming
  - Require the prior knowledge of the environment's dynamics, the complete probability distributions of all possible transitions
  - In many cases, it is easy to generate experience sampled according to the desired probability distributions, but **infeasible to obtain the distributions in explicit form**
- Model-free algorithms, e.g., Monte Carlo methods, temporal-difference learning
  - Require only **experience** – sample sequences of states, actions, and rewards from actual or simulated interaction with an environment
  - **Learning** from actual experience is striking because it requires no prior knowledge of the environment's dynamics, yet can still attain optimal behavior

# The term “Monte Carlo” (MC)

- A broad class of computational algorithms that rely on repeated random sampling to obtain numerical results
  - Use randomness to solve problems that might be deterministic in principle
- The **Law of Large Numbers**: the basis for Monte Carlo simulations
  - As the number of identically distributed, randomly generated variables increases, their sample mean approaches their theoretical mean

# Monte Carlo estimate of PI

- We can estimate pi to as many digits as we like by simply playing a game of darts
  - Generate random points within a box, and counting the number of points which fall within an embedded circle



# Monte Carlo methods for RL

- Based on averaging sample returns
  - We define Monte Carlo methods only for episodic tasks, to ensure that well-defined returns are available
  - Incremental in an **episode-by-episode** sense, but not in a step-by-step sense
  - Average complete returns, as opposed to methods that learn from partial returns, e.g., temporal-difference learning

# Monte Carlo prediction

- Considering Monte Carlo methods for learning the state-value function for a given policy
  - $V_\pi(s)$ : the expected return—expected cumulative future discounted reward—starting from  $s$
  - Estimate  $V_\pi(s)$  from **experience**: simply average the returns observed after visits to  $s$
  - As more returns are observed, the average should converge to the expected value

$$\begin{aligned}V_\pi(s) &= \mathbb{E}_\pi[G_t | S_t = s] \\&= \mathbb{E}_\pi[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]\end{aligned}$$

# Monte Carlo prediction

- Some notations
  - Each occurrence of state  $s$  in an episode is called a **visit** to  $s$
  - $s$  may be visited multiple times in the same episode, we call the first time it is visited in an episode the **first visit** to  $s$
- First-visit Monte Carlo method
  - Estimate  $V_\pi(s)$  as the average of the returns following first visits to  $s$
  - Have been most widely studied
- Every-visit Monte Carlo method
  - Average the returns following visits to  $s$
  - Extend more naturally to function approximation and eligibility traces

# First-visit MC prediction

## First-visit MC prediction, for estimating $V \approx v_\pi$

Input: a policy  $\pi$  to be evaluated

Initialize:

$V(s) \in \mathbb{R}$ , arbitrarily, for all  $s \in \mathcal{S}$

$Returns(s) \leftarrow$  an empty list, for all  $s \in \mathcal{S}$

Loop forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

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

Unless  $S_t$  appears in  $S_0, S_1, \dots, S_{t-1}$ :

Append  $G$  to  $Returns(S_t)$

$V(S_t) \leftarrow \text{average}(Returns(S_t))$

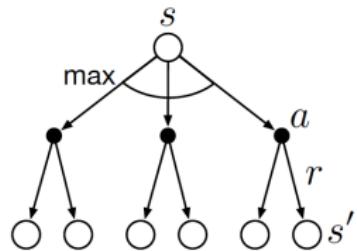
# First-visit MC prediction

- Both first-visit MC and every-visit MC converge to  $V_\pi$  as the number of (first) visits to  $s$  goes to infinity
- For first-visit MC
  - Each return is an independent, identically distributed estimate of  $V_\pi$  with finite variance
  - Each average is an unbiased estimate, and the standard deviation of its error falls as  $1/\sqrt{n}$

# Backup diagram for MC prediction



- The root is a state node, and below it is the entire trajectory of transitions along a particular single episode, ending at the terminal state
- MC shows only those sampled on the one episode
  - Opposed to DP that shows all possible transitions
- MC goes all the way to the end of the episode
  - Opposed to DP that includes only one-step transitions



# Properties of MC

- MC methods do not bootstrap
  - Estimates for each state are independent
  - The estimate for one state does not build upon the estimate of any other state
- The computational expense of estimating the value of a single state is independent of the number of states
  - This can make MC particularly attractive when one requires the value of only one or a subset of states
  - One can generate many sample episodes starting from the states of interest, averaging returns from only these states, ignoring all others

# MC estimation of action values

- With a model
  - State values alone are sufficient to determine a policy
- Without a model
  - One must explicitly estimate the value of each action in order for the values to be useful in suggesting a policy
- A state-action pair  $(s, a)$  is said to be visited in an episode if ever the state  $s$  is visited and action  $a$  is taken in it
  - To help in choosing among the actions available in each state, we need to estimate the value of all the actions from each state, not just the one we currently favor
  - This is the general problem of **maintaining exploration**

# Maintaining exploration

- The assumption of **exploring starts**
  - Specify that the episodes start in a state-action pair, and that every pair has **a nonzero probability** of being selected as the start
  - Guarantee that all state-action pairs will be visited **an infinite number of times** in the limit of an infinite number of episodes
  - Sometimes useful, but cannot be relied upon in general, particularly when learning directly from actual interaction with an environment
- The most common alternative approach to assuring that all state-action pairs are encountered is to consider only policies that are **stochastic** with a nonzero probability of selecting all actions in each state

# Table of Contents

## 1 Monte Carlo Methods

- Monte Carlo prediction
- Monte Carlo control

## 2 Temporal-Difference Learning

- Temporal-difference prediction
- Temporal-difference control

## 3 Extensions

- $n$ -step bootstrapping
- Eligibility traces

# Monte Carlo control

- Construct each  $\pi_{k+1}$  as the greedy policy w.r.t.  $Q_{\pi_k}$

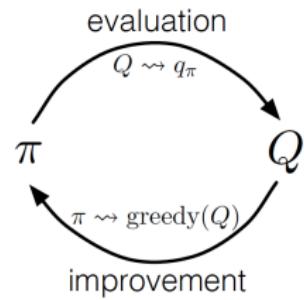
$$\pi_{k+1}(s) = \arg \max_a Q_{\pi_k}(s, a)$$

- Policy improvement theorem

$$\begin{aligned} Q_{\pi_k}(s, \pi_{k+1}(s)) &= Q_{\pi_k}(s, \arg \max_a Q_{\pi_k}(s, a)) = \max_a Q_{\pi_k}(s, a) \\ &\geq Q_{\pi_k}(s, \pi_k(s)) = V_{\pi_k}(s) \end{aligned}$$

- Use **generalized policy iteration** to approach optimal policies

$$\pi_0 \xrightarrow{E} Q_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} Q_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \dots \xrightarrow{I} \pi_* \xrightarrow{E} Q_*$$



# MC with policy iteration

Monte Carlo ES (Exploring Starts), for estimating  $\pi \approx \pi_*$

Initialize:

$\pi(s) \in \mathcal{A}(s)$  (arbitrarily), for all  $s \in \mathcal{S}$

$Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$  empty list, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Loop forever (for each episode):

Choose  $S_0 \in \mathcal{S}, A_0 \in \mathcal{A}(S_0)$  randomly such that all pairs have probability  $> 0$

Generate an episode from  $S_0, A_0$ , following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

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

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :

Append  $G$  to  $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$

- Two unlikely assumptions
  - The episodes have exploring starts
  - Policy evaluation could be done with an infinite number of episodes

# MC control without exploring starts

- Without the assumption of exploring starts
  - If simply making the policy greedy w.r.t. the current value function, it can meet the policy improvement theorem, but would **prevent further exploration** of non-greedy actions
- Use a **soft** policy
  - $\pi(a|s) > 0, \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$
  - Gradually shifted closer and closer to a deterministic optimal policy
  - $\varepsilon$ -soft policies:  $\pi(a|s) \geq \frac{\varepsilon}{|\mathcal{A}(s)|}$
- $\varepsilon$ -greedy policy
  - With probability  $\varepsilon$ , select an action at random
  - Otherwise, choose the greedy action
  - An example of the  $\varepsilon$ -soft policy, closest to greedy among  $\varepsilon$ -soft policies

# MC control with soft policies

On-policy first-visit MC control (for  $\varepsilon$ -soft policies), estimates  $\pi \approx \pi_*$

Algorithm parameter: small  $\varepsilon > 0$

Initialize:

$\pi \leftarrow$  an arbitrary  $\varepsilon$ -soft policy

$Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$  empty list, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

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

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :

Append  $G$  to  $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken arbitrarily)

For all  $a \in \mathcal{A}(S_t)$ :

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

# Policy improvement theorem for soft policies

- For any  $\varepsilon$ -soft policy,  $\pi$ , any  $\epsilon$ -greedy policy,  $\pi'$ , w.r.t.  $Q_\pi$  is guaranteed to be better than or equal to  $\pi$

$$\begin{aligned} Q_\pi(s, \pi'(s)) &= \sum_a \pi'(a|s) Q_\pi(s, a) \\ &= \frac{\varepsilon}{|\mathcal{A}(s)|} \sum_a Q_\pi(s, a) + (1 - \varepsilon) \max_a Q_\pi(s, a) \\ &\geq \frac{\varepsilon}{|\mathcal{A}(s)|} \sum_a Q_\pi(s, a) + (1 - \varepsilon) \sum_a \frac{\pi(a|s) - \frac{\varepsilon}{|\mathcal{A}(s)|}}{1 - \varepsilon} Q_\pi(s, a) \\ &= \frac{\varepsilon}{|\mathcal{A}(s)|} \sum_a Q_\pi(s, a) - \frac{\varepsilon}{|\mathcal{A}(s)|} \sum_a Q_\pi(s, a) + \sum_a \pi(a|s) Q_\pi(s, a) \\ &= v_\pi(s) \end{aligned}$$

## ARTICLE

doi:10.1038/nature16961

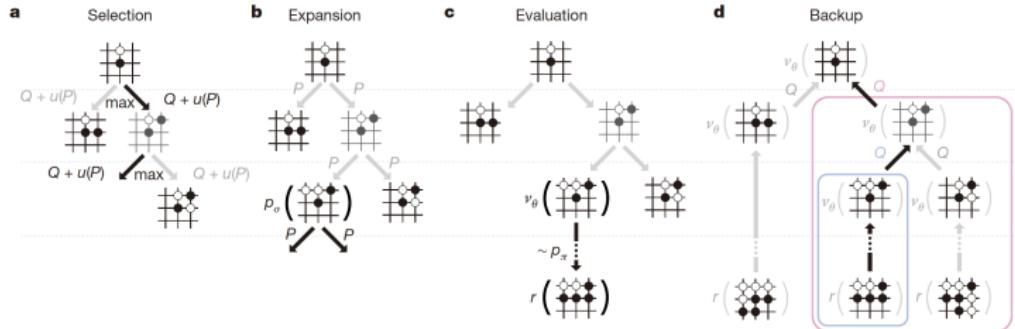
# Mastering the game of Go with deep neural networks and tree search

David Silver<sup>1\*</sup>, Aja Huang<sup>1\*</sup>, Chris J. Maddison<sup>1</sup>, Arthur Guez<sup>1</sup>, Laurent Sifre<sup>1</sup>, George van den Driessche<sup>1</sup>, Julian Schrittwieser<sup>1</sup>, Ioannis Antonoglou<sup>1</sup>, Veda Panneershelvam<sup>1</sup>, Marc Lanctot<sup>1</sup>, Sander Dieleman<sup>1</sup>, Dominik Grewe<sup>1</sup>, John Nham<sup>2</sup>, Nal Kalchbrenner<sup>1</sup>, Ilya Sutskever<sup>2</sup>, Timothy Lillicrap<sup>1</sup>, Madeleine Leach<sup>1</sup>, Koray Kavukcuoglu<sup>1</sup>, Thore Graepel<sup>1</sup> & Demis Hassabis<sup>1</sup>

Monte Carlo tree search (MCTS)<sup>11,12</sup> uses Monte Carlo rollouts to estimate the value of each state in a search tree. As more simulations are executed, the search tree grows larger and the relevant values become more accurate. The policy used to select actions during search is also improved over time, by selecting children with higher values. Asymptotically, this policy converges to optimal play, and the evaluations converge to the optimal value function<sup>12</sup>. The strongest current Go programs are based on MCTS, enhanced by policies that are trained to predict human expert moves<sup>13</sup>. These policies are used

# MC in AlphaGo

RESEARCH ARTICLE



**Figure 3 | Monte Carlo tree search in AlphaGo.** a, Each simulation traverses the tree by selecting the edge with maximum action value  $Q$ , plus a bonus  $u(P)$  that depends on a stored prior probability  $P$  for that edge. b, The leaf node may be expanded; the new node is processed once by the policy network  $p_\sigma$  and the output probabilities are stored as prior probabilities  $P$  for each action. c, At the end of a simulation, the leaf node

is evaluated in two ways: using the value network  $v_\theta$ ; and by running a rollout to the end of the game with the fast rollout policy  $p_\pi$ , then computing the winner with function  $r$ . d, Action values  $Q$  are updated to track the mean value of all evaluations  $r(\cdot)$  and  $v_\theta(\cdot)$  in the subtree below that action.

- Instead of brute forcing from millions of possible ways to find the right path, **Monte Carlo Tree Search** algorithm chooses the best possible move from the current state of the game's tree with the help of RL.

# Review

- Monte Carlo methods
  - Model-free, directly learn from experience
  - Do not bootstrap
- First-visit MC prediction
  - Should maintain exploration for estimation of action values
- MC control with exploring starts
  - Can use a soft policy, i.e.,  $\varepsilon$ -greedy policy
  - Meet the policy improvement theorem for  $\varepsilon$ -soft policy



# Table of Contents

## 1 Monte Carlo Methods

- Monte Carlo prediction
- Monte Carlo control

## 2 Temporal-Difference Learning

- Temporal-difference prediction
- Temporal-difference control

## 3 Extensions

- $n$ -step bootstrapping
- Eligibility traces

# Temporal-Difference (TD) Learning

- If one had to identify one idea as **central** and novel to RL, it would undoubtedly be temporal-difference learning
  - A **combination** of Monte Carlo ideas and dynamic programming ideas
  - Like MC, TD can learn directly from raw experience without a model of the environment's dynamics (**model-free**)
  - Like DP, TD updates estimates based in part on other learned estimates, without waiting for a final outcome (**bootstrap**)

# TD prediction

- MC and TD in common
  - Use **experience** to solve the prediction problem, update their estimate  $V$  of  $V_\pi$  for the non-terminal state  $S_t$  occurring in that experience
- MC: must wait until the return following the visit is known (end of an episode), then use that return as a target for  $V(S_t)$

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

- TD: need to wait only until the next time step, use  $R_{t+1} + \gamma V(S_{t+1})$  as the target for  $V(S_t)$ , bootstrapping

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

# Basic TD algorithm

## Tabular TD(0) for estimating $v_\pi$

Input: the policy  $\pi$  to be evaluated

Algorithm parameter: step size  $\alpha \in (0, 1]$

Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$

Loop for each episode:

    Initialize  $S$

    Loop for each step of episode:

$A \leftarrow$  action given by  $\pi$  for  $S$

        Take action  $A$ , observe  $R, S'$

$V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

    until  $S$  is terminal

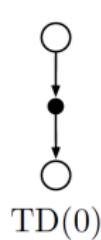
- **TD(0):** A special case of the TD( $\lambda$ ) method using eligibility trace
- **One-step TD:** A special case of the  $n$ -step TD methods

# Deep insight of the Bellman equation

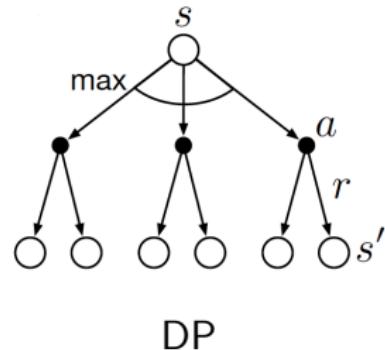
$$\begin{aligned}V_{\pi}(s) &= \mathbb{E}_{\pi}[G_t | S_t = s] \\&= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s] \\&= \mathbb{E}_{\pi}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s] \\&= \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V_{\pi}(s')]\end{aligned}$$

- **MC**: The expected  $G_t$  is not known, a sample return is used in place of the real expected return
- **DP**: The true  $V_{\pi}$  is not known, and the current estimate  $V(S_{t+1})$  is used instead
- **TD**: It samples the expected values  $R_{t+1}$ , and it uses the current estimate  $V(S_{t+1})$  instead of the true  $V_{\pi}$ 
  - Combine the sampling of MC with the bootstrapping of DP

# Backup diagram for TD(0)



- The value estimate for the state node is updated on the basis of the one sample transition from it to the immediately following state
- **Sample updates:** based on a single sample successor, involve looking ahead to a sample successor state



- **Expected updates:** based on a complete distribution of all possible successors

## TD error

- The difference between the estimated value of  $S_t$  and the better estimate  $R_{t+1} + \gamma V(S_{t+1})$

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

- The MC error, with  $V$  not updated during the episode

$$\begin{aligned} G_t - V(S_t) &= R_{t+1} + \gamma G_{t+1} - V(S_t) + \gamma V(S_{t+1}) - \gamma V(S_{t+1}) \\ &= \delta_t + \gamma(G_{t+1} - V(S_{t+1})) \\ &= \delta_t + \gamma \delta_{t+1} + \gamma^2(G_{t+2} - V(S_{t+2})) \\ &= \delta_t + \gamma \delta_{t+1} + \dots + \gamma^{T-t-1} \delta_{T-1} + \gamma^{T-t}(G_T - V(S_T)) \\ &= \delta_t + \gamma \delta_{t+1} + \dots + \gamma^{T-t-1} \delta_{T-1} + \gamma^{T-t}(0 - 0) \\ &= \sum_{k=t}^{T-1} \gamma^{k-t} \delta_k \end{aligned}$$

# Advantages of TD prediction methods

- Compared to DP, **do not require a model of the environment**, of its reward and next-state probability distribution
  - In many cases, it is infeasible to obtain the distributions in explicit form
- Compared to MC, **wait only one time step** instead of waiting until the end of an episode
  - Applications have very long episodes, or continuing tasks have no episodes at all

# Convergence property of TD prediction

- Are TD methods sound? Convergence guarantee?
  - It is convenient to learn one guess from the next, without waiting for an actual outcome
  - We can still **guarantee convergence** to the correct answer
  - For any fixed policy  $\pi$ , TD(0) has been proved to converge to  $V_\pi$ , in the mean for a constant step-size parameter if it is sufficiently small, and with probability 1 if the step-size parameter decreases according to the usual stochastic approximation conditions

# Convergence guarantee of TD prediction

- It is convenient to learn one guess from the next, without waiting for an actual outcome, and we can still guarantee convergence to the correct answer
  - For any fixed policy  $\pi$ , TD(0) has been proved to converge to  $V_\pi$ , in the mean for a constant step-size parameter if it is sufficiently small, and with probability 1 if the step-size parameter decreases according to the usual stochastic approximation conditions

- 
- Use the **Bellman equation** for  $v_\pi$  as an update rule

$$\begin{aligned}v_{k+1}(s) &= \mathbb{E}_\pi[R_{t+1} + \gamma v_k(S_{t+1})|S_t = s] \\&= \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a)[r + \gamma v_k(s')]\end{aligned}$$

- $v_k = v_\pi$  is a fixed point for this update rule
  - The sequence  $\{v_k\}$  converges to  $v_\pi$  as  $k \rightarrow \infty$  under the same conditions that guarantee the existence of  $v_\pi$

# Table of Contents

## 1 Monte Carlo Methods

- Monte Carlo prediction
- Monte Carlo control

## 2 Temporal-Difference Learning

- Temporal-difference prediction
- Temporal-difference control

## 3 Extensions

- $n$ -step bootstrapping
- Eligibility traces

# On-policy vs. off-policy

Target policy $\pi(a s)$	Behavior policy $b(a s)$
To be evaluated or improved	To explore to generate data
Make decisions finally	Make decisions in training phase

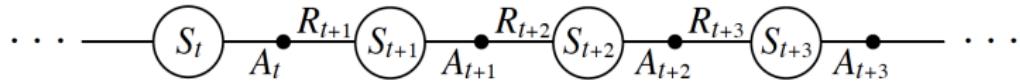
- **On-policy** methods:  $\pi(a|s) = b(a|s)$ 
  - Evaluate or improve the policy that is used to make decisions during training
  - e.g., SARSA
- **Off-policy** methods:  $\pi(a|s) \neq b(a|s)$ 
  - Evaluate or improve a policy different from that used to generate the data
  - Separate exploration from control
  - e.g., Q-learning

# SARSA: On-policy TD control

- We follow the pattern of generalized policy iteration
  - Only this time using TD methods for the evaluation or prediction part
- Learn the values of state-action pairs instead of a state-value function
  - Consider transitions between state-action pairs
- For an on-policy method, estimate  $Q_\pi(s, a)$  for the current behavior policy  $\pi$  and for all state-action pairs

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

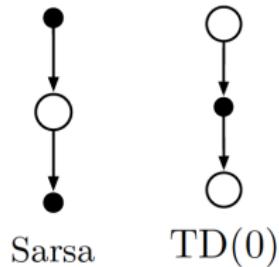
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$



# SARSA: On-policy TD control

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

- This update is done after every transition from non-terminal state  $S_t$ 
  - If  $S_{t+1}$  is terminal, then  $Q(S_{t+1}, A_{t+1})$  is defined as zero
  - Use every element of the quintuple of events,  $(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$



- The theorems assuring the convergence of state values under TD(0) also apply to the corresponding algorithm for action values

# SARSA algorithm

Sarsa (on-policy TD control) for estimating  $Q \approx q_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

    Initialize  $S$

    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

    Loop for each step of episode:

        Take action  $A$ , observe  $R, S'$

        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$$

$S \leftarrow S'; A \leftarrow A'$ ;

    until  $S$  is terminal

- Follow the pattern of generalized policy iteration
  - Continually estimate  $Q_\pi$  for the behavior policy  $\pi$ , and at the same time change  $\pi$  toward greediness w.r.t.  $Q_\pi$

# Convergence properties of SARSA

- Let behavior policy  $\pi$  be greedy w.r.t. current action value function
  - Meet the policy improvement theorem
  - However, prevent further exploration of non-greedy actions
- Instead, **soft** policies are favored
  - $\varepsilon$ -greedy or  $\varepsilon$ -soft policies
- Convergence properties of SARSA depend on the nature of the policy's dependence on the action-value function  $Q$ 
  - SARSA converges with probability 1 to an optimal policy and action-value function as long as all state-action pairs are visited an infinite number of times and the policy converges in the limit to the greedy policy (which can be arranged, for example, with  $\varepsilon$ -greedy policies by setting  $\varepsilon = 1/t$ )

# Compare the convergence property of DP

- Use the **Bellman equation** for  $v_\pi$  as an update rule

$$\begin{aligned}v_{k+1}(s) &= \mathbb{E}_\pi[R_{t+1} + \gamma v_k(S_{t+1})|S_t = s] \\&= \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v_k(s')]\end{aligned}$$

- $v_k = v_\pi$  is a fixed point for this update rule
  - The sequence  $\{v_k\}$  converges to  $v_\pi$  as  $k \rightarrow \infty$  under the same conditions that guarantee the existence of  $v_\pi$

- 
- Using policy improvement theorem, we can obtain a sequence of monotonically improving policies and value functions

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \dots \xrightarrow{I} \pi_* \xrightarrow{E} v_*$$

- This process is guaranteed to converge to an optimal policy and optimal value function in a finite number of iterations
  - Each policy is guaranteed to be a strictly improvement over the previous one unless it is already optimal
  - A finite MDP has only a finite number of policies

# Q-learning: Off-policy TD control

- One of the early breakthroughs in RL was the development of an off-policy TD control algorithm, known as **Q-learning**

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

- The learned action-value function,  $Q$ , directly **approximates the optimal action-value function**,  $Q_*$ , independent of the policy being followed

# Q-learning: Off-policy TD control

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

- Use the Bellman optimality equation into an update rule

$$Q_*(S_t, A_t) = R_{t+1} + \gamma \max_a Q_*(S_{t+1}, a)$$

- 
- Value iteration

$$v_{k+1}(s) = \max_a \sum_{s',r} p(s', r | s, a) [r + \gamma v_k(s')]$$

- Turn Bellman optimality equation into an update rule
- Directly approximate the optimal state-value function,  $v_*$

# Q-learning algorithm

Q-learning (off-policy TD control) for estimating  $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

    Initialize  $S$

    Loop for each step of episode:

        Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

        Take action  $A$ , observe  $R, S'$

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

$S \leftarrow S'$

    until  $S$  is terminal

# The reason for “off-policy”

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

- Q-learning approximates the optimal action-value function for an optimal policy,  $Q \approx Q_* = Q_{\pi_*}$ 
    - The target policy is greedy w.r.t  $Q$ ,  $\pi(a|s) = \arg \max_a Q(s, a)$
    - The behavior policy can be others, e.g.,  $b(a|s) = \varepsilon$ -greedy
- 

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

- SARSA approximates the action-value function for the behavior policy,  $Q \approx Q_\pi = Q_b$ 
  - The target and the behavior policy are the same, e.g.,  $\pi(a|s) = b(a|s) = \varepsilon$ -greedy

# Convergence properties of Q-learning

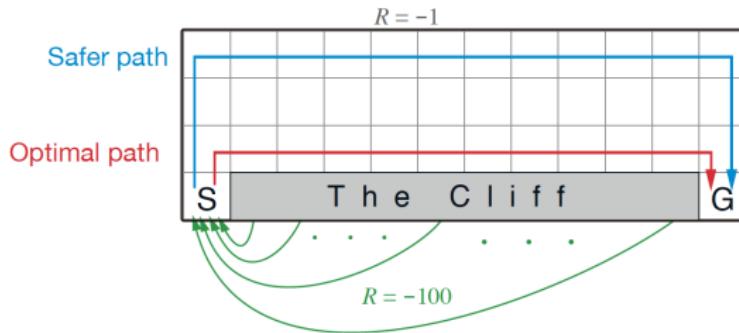
- All state-action pairs are visited an infinite number of times
  - This is a minimal requirement in the sense that any method guaranteed to find optimal behavior in the general case must require it
  - Under this assumption and a variant of the usual stochastic approximation conditions on the sequence of step-size parameters,  $Q$  has been shown to converge with probability 1 to  $Q_*$

- 
- Use the **Bellman equation** for  $v_\pi$  as an update rule

$$\begin{aligned}v_{k+1}(s) &= \mathbb{E}_\pi[R_{t+1} + \gamma v_k(S_{t+1})|S_t = s] \\&= \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a)[r + \gamma v_k(s')]\end{aligned}$$

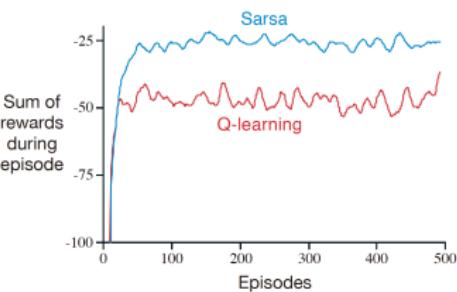
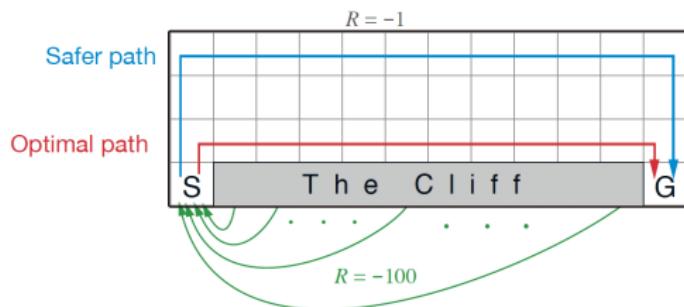
- $v_k = v_\pi$  is a fixed point for this update rule
  - The sequence  $\{v_k\}$  converges to  $v_\pi$  as  $k \rightarrow \infty$  under the same conditions that guarantee the existence of  $v_\pi$

## Example: Cliff walking



- The usual action causing movement up, down, right, and left
- $r = -1$  on all transitions except those into the region marked “The Cliff”. Stepping into this region incurs a reward of  $-100$  and sends the agent instantly back to the start.
- With  $\varepsilon$ -greedy action selection,  $\varepsilon = 0.1$

# On-policy vs. Off-policy learning



- Q-learning learns values for the optimal policy, traveling right along the edge of the cliff
  - Occasionally fall off the cliff because of the  $\varepsilon$ -greedy action selection
- SARSA takes the action selection into account and learns the longer but safer path through the upper part of the grid
  - Q-learning's online performance is worse than that of SARSA
  - If  $\varepsilon$  is gradually reduced, both methods would asymptotically converge to the optimal policy

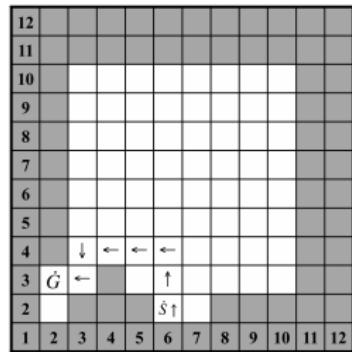
# Code example for maze navigation

```
class QLearning(object):
    def __init__(self, env, gamma=0.99, lr=0.1):
        super(QLearning, self).__init__()
        self.gamma = gamma; self.lr = lr; self.env = env
        self.Q = np.zeros((env.observation_space.n, env.action_space.n))

    def greedy(self, s, Q):
        candi_actions = []
        Qmax = -np.inf
        for a in range(self.env.action_space.n):
            if Q[s, a] > Qmax:
                candi_actions = []; candi_actions.append(a)
                Qmax = Q[s, a]
            elif Q[s, a] == Qmax:
                candi_actions.append(a)
        return np.random.choice(candi_actions)

    # epsilon-greedy policy
    def pi(self, s, epsilon=0.1):
        if np.random.random() < epsilon:
            ### random action
            return self.env.action_space.sample()
        else:
            ### epsilon greedy
            return self.greedy(s, self.Q)

    def step(self, s, a, r, s_next):
        Qmax_s_next = np.max(self.Q[s_next])
        self.Q[s, a] += self.lr * (r + self.gamma * Qmax_s_next - self.Q[s, a])
```



# Code example for maze navigation

```
learner = QLearning(env, gamma=gamma, lr=lr)

# record the navigation steps in each episode
steps = np.zeros(max_epochs)

for epoch in tqdm(range(max_epochs)):
    s = env.reset()
    for step in range(max_steps):

        # selection action according to epsilon-greedy policy
        a = learner.pi(s, epsilon=epsilon)

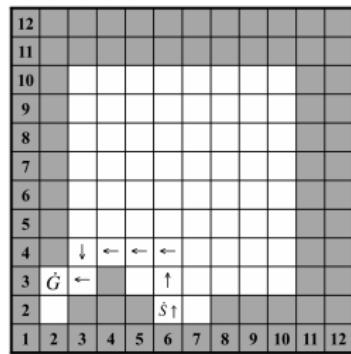
        # interact with the environment
        s_next, r, done, _ = env.step(a)

        # update Q-function
        learner.step(s, a, r, s_next)

        # whether transit to the terminal state
        if done: break

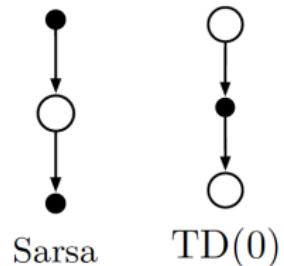
        # transit to next state
        s = s_next

    steps[epoch] = step
```



# Review

- TD learning
  - A combination of DP (bootstrap) and MC (model free)
  - Guarantee convergence under certain assumptions
- On-policy vs. off-policy control
  - SARSA: on-policy, behavior policy = target policy
  - Q-learning: off-policy, behavior policy  $\neq$  target policy



# Table of Contents

## 1 Monte Carlo Methods

- Monte Carlo prediction
- Monte Carlo control

## 2 Temporal-Difference Learning

- Temporal-difference prediction
- Temporal-difference control

## 3 Extensions

- $n$ -step bootstrapping
- Eligibility traces

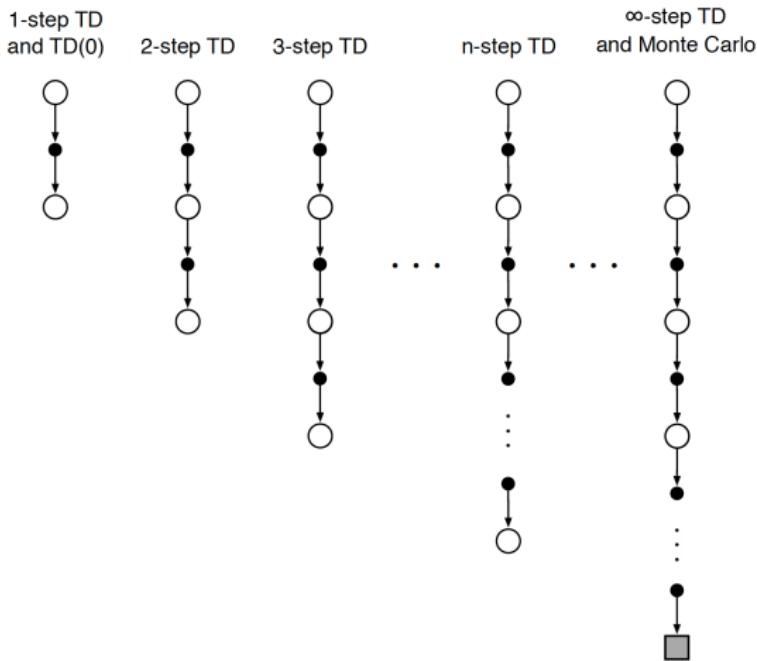
## Combine MC and one-step TD

- Neither MC or one-step TD is always the best, we **generalize** both methods so that one can shift from one to the other smoothly as needed to meet the demands of a particular task
- One-step TD: In many applications, one wants to be able to update the action very fast to take into account anything that has changed
- However, bootstrapping works best if it is over a length of time in which a significant and recognizable state change has occurred

$n = 1$	$n$ -step TD	$n = \infty$
$\text{TD}(0)$	$\leftrightarrow$	$MC$

# $n$ -step TD prediction

- Perform an update based on an intermediate number of rewards, more than one, but less than all of them until termination



## Recall MC and TD(0) updates

- In MC updates, the target is the **complete return**

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

$$\begin{aligned} V(S_t) &\leftarrow V(S_t) + \alpha[\textcolor{red}{G_t} - V(S_t)] \\ &= V(S_t) + \alpha[\textcolor{red}{R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t+1} R_T} - V(S_t)] \end{aligned}$$

- In TD(0) updates, the target is the **one-step return**

$$G_{t:t+1} = R_{t+1} + \gamma V(S_{t+1})$$

$$\begin{aligned} V(S_t) &\leftarrow V(S_t) + \alpha[\textcolor{red}{G_{t:t+1}} - V(S_t)] \\ &= V(S_t) + \alpha[\textcolor{red}{R_{t+1} + \gamma V(S_{t+1})} - V(S_t)] \end{aligned}$$

## $n$ -step TD update rule

- For  $n$ -step TD, set the target as the  **$n$ -step return**

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

- All  $n$ -step returns can be considered approximations to the complete return, truncated after  $n$  steps and then corrected for the remaining missing terms by  $V(S_{t+n})$

$$V(S_t) \leftarrow V(S_t) + \alpha[G_{t:t+n} - V(S_t)]$$

$$= V(S_t) + \alpha[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n}) - V(S_t)]$$

## *n*-step TD algorithm

*n*-step TD for estimating  $V \approx v_\pi$

Input: a policy  $\pi$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , a positive integer  $n$

Initialize  $V(s)$  arbitrarily, for all  $s \in \mathcal{S}$

All store and access operations (for  $S_t$  and  $R_t$ ) can take their index mod  $n + 1$

Loop for each episode:

Initialize and store  $S_0 \neq \text{terminal}$

$$T \leftarrow \infty$$

Loop for  $t = 0, 1, 2, \dots$ :

| If  $t < T$ , then:

Take an action according to  $\pi(\cdot | S_t)$

Observe and store the next reward as  $R_{t+1}$  and the next state as  $S_{t+1}$

If  $S_{t+1}$  is terminal, then  $T \leftarrow t + 1$

$\tau \leftarrow t - n + 1$     ( $\tau$  is the time whose state's estimate is being updated)

If  $\tau \geq 0$ :

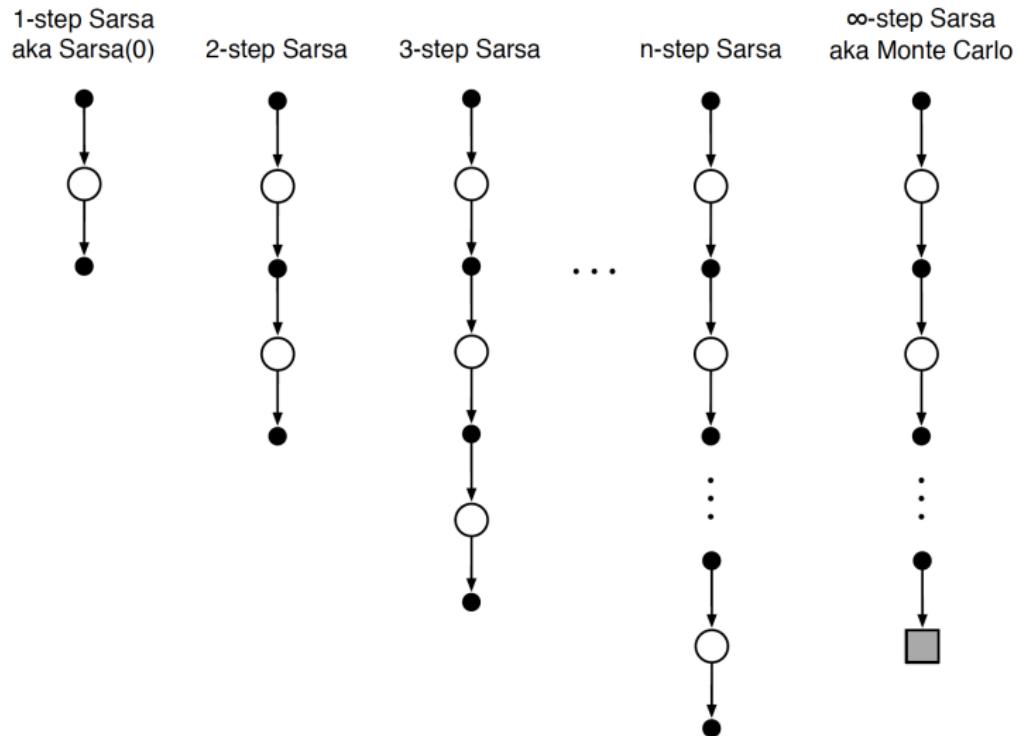
$$G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$$

If  $\tau + n < T$ , then:  $G \leftarrow G + \gamma^n V(S_{\tau+n})$

$$V(S_\tau) \leftarrow V(S_\tau) + \alpha [G - V(S_\tau)]$$

Until  $\tau = T - 1$

# $n$ -step TD for on-policy control: $n$ -step SARSA



# From $n$ -step TD prediction to $n$ -step SARSA

- The main idea is to simply switch states for state-action pairs, and then use an  $\varepsilon$ -greedy policy
- re-define  $n$ -step returns (update targets) in terms of estimated action values

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n}, A_{t+n})$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [G_{t:t+n} - Q(S_t, A_t)]$$

$$= Q(S_t, A_t) + \alpha [R_{t+1} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n}, A_{t+n}) - Q(S_t, A_t)]$$

## *n*-step SARSA algorithm

*n*-step Sarsa for estimating  $Q \approx q_*$  or  $q_\pi$

Initialize  $Q(s, a)$  arbitrarily, for all  $s \in \mathcal{S}, a \in \mathcal{A}$

Initialize  $\pi$  to be  $\varepsilon$ -greedy with respect to  $Q$ , or to a fixed given policy

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$ , a positive integer  $n$

All store and access operations (for  $S_t$ ,  $A_t$ , and  $R_t$ ) can take their index mod  $n + 1$ .

Loop for each episode:

Initialize and store  $S_0 \neq \text{terminal}$

Select and store an action  $A_0 \sim \pi(\cdot | S_0)$

$$T \leftarrow \infty$$

Loop for  $t \equiv 0, 1, 2, \dots$ :

| If  $t < T$ , then:

Take action. At

Observe and store the next reward as  $R_{t+1}$  and the next state as  $S_{t+1}$

If  $S_{t+1}$  is terminal, then:

$$T \leftarrow t + 1$$

else:

Select and store an action  $A_{t+1} \sim \pi(\cdot | S_{t+1})$

$\tau \leftarrow t - n + 1$     ( $\tau$  is the time whose estimate is being updated)

If  $\tau > 0$ :

$$\bar{G} \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$$

If  $\tau + n < T$ , then  $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$

$$(G_{\tau:\tau+n})$$

$$Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha [G - Q(S_\tau, A_\tau)]$$

If  $\pi$  is being learned, then ensure that  $\pi(\cdot|S_\tau)$  is  $\varepsilon$ -greedy wrt  $Q$

Until  $\tau = T - 1$

# Table of Contents

## 1 Monte Carlo Methods

- Monte Carlo prediction
- Monte Carlo control

## 2 Temporal-Difference Learning

- Temporal-difference prediction
- Temporal-difference control

## 3 Extensions

- $n$ -step bootstrapping
- Eligibility traces

# Eligibility traces: unify/generalize TD and MC

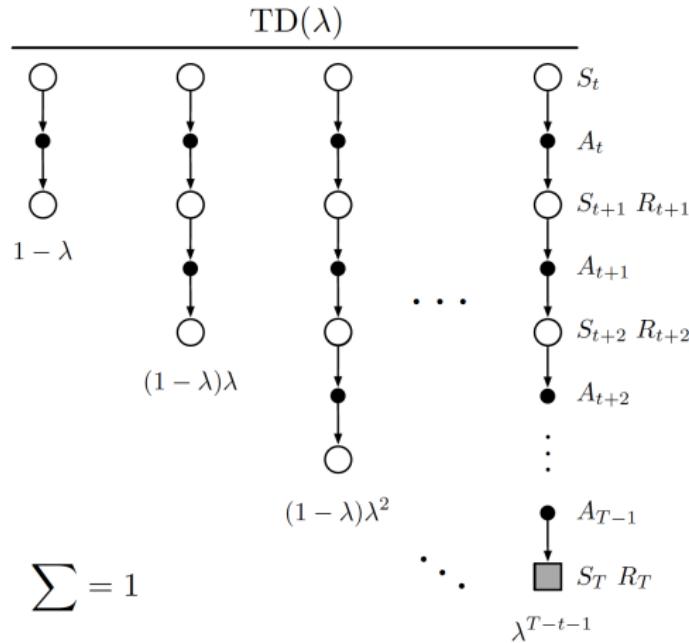
- Almost any TD method can be combined with eligibility traces to obtain a more general method that may learn more efficiently
  - e.g., the popular  $\text{TD}(\lambda)$  algorithm,  $\lambda$  refers the use of an eligibility trace
  - Produce a family of methods spanning a spectrum that has MC methods at one end ( $\lambda = 1$ ) and one-step TD methods at the other ( $\lambda = 0$ )
- Eligibility traces offer an elegant algorithmic mechanism with significant computational advantages (compared to  $n$ -step TD)
  - Only a single trace vector is required rather than a store of the last  $n$  feature vectors
  - Learning also occurs continually and uniformly in time rather than being delayed and then catching up at the end of the episode
  - Learning can occur and effect behavior immediately after a state is encountered rather than being delayed  $n$ -steps

# The $\lambda$ -return

- How to interrelate TD and MC?
  - e.g., average one-step and infinite-step returns,  $G = (G_t + G_{t:t+1})/2$
  - An update that averages simpler component updates is called a **compound update**
- The TD( $\lambda$ ) algorithm can be understood as one particular way of averaging  $n$ -step updates

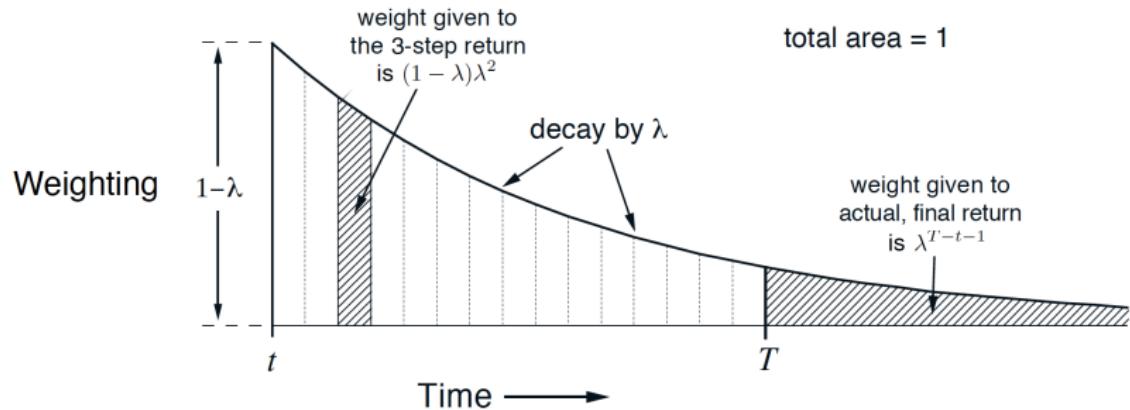
$$\begin{aligned} G_t^\lambda &= (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{t:t+n} \\ &= (1 - \lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_{t:t+n} + \lambda^{T-t-1} G_t \end{aligned}$$

# Backup diagram for TD( $\lambda$ )



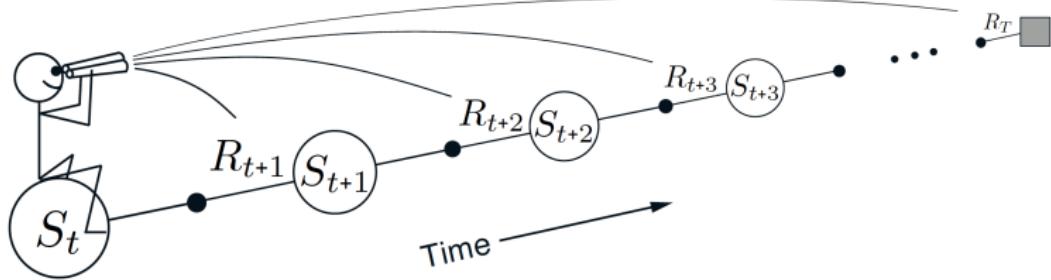
**Figure 12.1:** The backup diagram for TD( $\lambda$ ). If  $\lambda = 0$ , then the overall update reduces to its first component, the one-step TD update, whereas if  $\lambda = 1$ , then the overall update reduces to its last component, the Monte Carlo update.

# The weight distribution



**Figure 12.2:** Weighting given in the  $\lambda$ -return to each of the  $n$ -step returns.

# The forward view



**Figure 12.4:** The forward view. We decide how to update each state by looking forward to future rewards and states.

## TD( $\lambda$ ): distribute the computation in each time step

- The eligibility trace,  $z(s), \forall s \in \mathcal{S}$ , in each episode:

$$z_0(s) \leftarrow 0, \quad \forall s \in \mathcal{S}$$

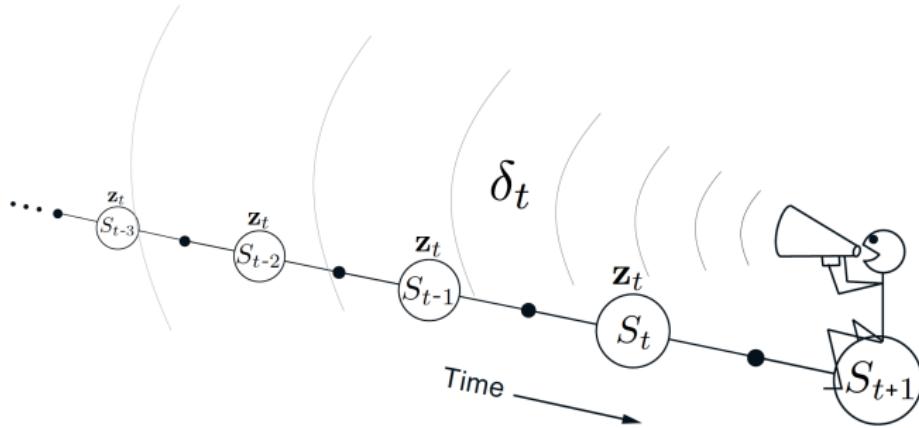
$$z_t(s) \leftarrow \gamma \lambda z_{t-1}(s) + \mathbb{I}(S_t = s), \quad \forall s \in \mathcal{S}$$

- The update rule for state-value function

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

$$V(s) \leftarrow V(s) + \alpha \delta_t \cdot \textcolor{red}{z_t(s)}, \quad \forall s \in \mathcal{S}$$

# The backward view



**Figure 12.5:** The backward or mechanistic view of  $\text{TD}(\lambda)$ . Each update depends on the current TD error combined with the current eligibility traces of past events.

- True online  $\text{TD}(\lambda)$  algorithms...
- $\text{TD}(\lambda)$  for control: SARSA( $\lambda$ ), Q( $\lambda$ )

# Learning objectives of this lecture

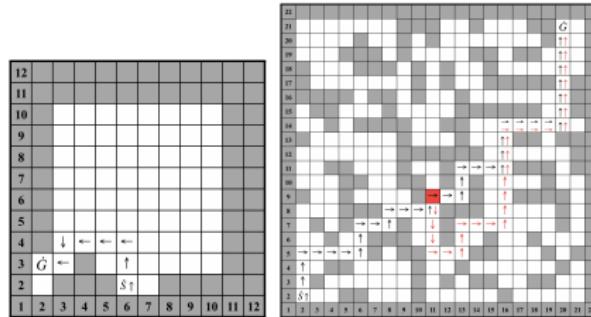
- You should be able to...
  - Understand and be able to use the typical three RL algorithms: DP, MC, and TD
  - Understand relationships and differences among the three algorithms
  - Know concepts of bootstrapping
  - Know concepts of model-free vs. model-based
  - Know concepts of on-policy vs. off-policy
  - Know concepts of prediction vs. control
  - Be aware of advanced TD methods, e.g.,  $n$ -step TD,  $\text{TD}(\lambda)$

# References

- Chapters 5 & 6 & 7 & 12, *Reinforcement Learning: An Introduction*, 2nd Edition
  - <http://202.119.32.195/cache/2/03/incompleteideas.net/5b682cef88335157bc5542267cf3f442/RLbook2018.pdf>
- Papers
  - C. J. C. H. Watkins and P. Dayan, “**Q-learning**,” *Machine Learning*, 1992.

# Homework 2

- Study the Q-learning algorithm in detail
- Implement the Q-learning algorithm on the following two maze navigation problems
  - The gray is obstacles,  $S$  and  $G$  are start and goal, respectively
  - Use **Python** as the programming language
- Write a report that introduces Q-learning and your experimentation
  - Explanations, steps, evaluation results, visualizations...
  - Submit the code and the report to [mg20150005@mail.nju.edu.cn](mailto:mg20150005@mail.nju.edu.cn)



# THE END