

Lecture 7: Review on Tabular RL

Bolei Zhou

The Chinese University of Hong Kong

bzhou@ie.cuhk.edu.hk

February 18, 2020

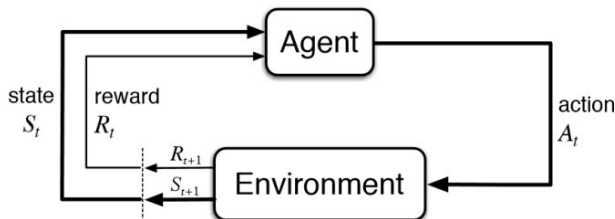
Announcement

- 1 We continue with ZOOM Online Teaching
 - 1 20-30 mins per session
 - 2 Try to open your webcam
 - 3 Take note
- 2 HW1 and project proposal due by the end of Feb.19, please email to cuhkrlcourse@googlegroups.com

Today's Content: Review Session

- ① RL Basics
- ② MDP prediction and control
 - ① Policy evaluation
 - ② Control: Policy iteration and value iteration
- ③ Model-free prediction and control

Reinforcement Learning



- 1 a computational approach to learning whereby **an agent** tries to **maximize** the total amount of **reward** it receives while interacting with a complex and uncertain **environment**.
 - Sutton and Barto

Difference between Reinforcement Learning and Supervised Learning

- ① Sequential data as input (not i.i.d)
- ② Trial-and-error exploration (balance between exploration and exploitation)
- ③ There is no supervisor, only a reward signal, which is also delayed
- ④ Agent's actions affect the subsequent data it receives (agent's action changes the environment)

Major Components of RL

- ① An RL agent may include one or more of these components
 - ① Policy: agent's behavior function
 - ② Value function: how good is each state or action
 - ③ Model: agent's state representation of the environment

Define the model of environment

- ① Markov Processes
- ② Markov Reward Processes(MRPs)
- ③ Markov Decision Processes (MDPs)

Markov Process

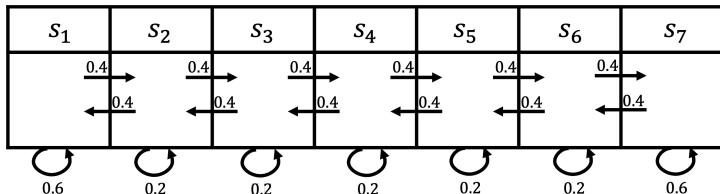
- ① P is the dynamics/transition model that specifies $p(s_{t+1} = s' | s_t = s)$
- ② State transition matrix: $P(\text{To}|\text{From})$

$$P = \begin{bmatrix} P(s_1|s_1) & P(s_2|s_1) & \dots & P(s_N|s_1) \\ P(s_1|s_2) & P(s_2|s_2) & \dots & P(s_N|s_2) \\ \vdots & \vdots & \ddots & \vdots \\ P(s_1|s_N) & P(s_2|s_N) & \dots & P(s_N|s_N) \end{bmatrix}$$

Markov Reward Process (MRP)

- ① Markov Reward Process is a Markov Chain + reward
- ② Definition of Markov Reward Process (MRP)
 - ① S is a (finite) set of states ($s \in S$)
 - ② P is dynamics/transition model that specifies $P(S_{t+1} = s' | s_t = s)$
 - ③ **R is a reward function** $R(s_t = s) = \mathbb{E}[r_t | s_t = s]$
 - ④ Discount factor $\gamma \in [0, 1]$
- ③ If finite number of states, R can be a vector

MRP Example



Reward: $+5$ in s_1 , $+10$ in s_7 , 0 in all other states. So that we can represent $R = [5, 0, 0, 0, 0, 0, 10]$

Return Function and Value Function

- ① Return: Discounted sum of rewards from time step t to horizon

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots + \gamma^{T-t-1} R_T$$

- ② Value function $V_t(s)$: Expected return from t in state s

$$\begin{aligned} V_t(s) &= \mathbb{E}[G_t | s_t = s] \\ &= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T | s_t = s] \end{aligned}$$

Computing the Value of a Markov Reward Process

- ① Value function: expected return from starting in state s

$$V(s) = \mathbb{E}[G_t | s_t = s] = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | s_t = s]$$

- ② MRP value function satisfies the following Bellman equation:

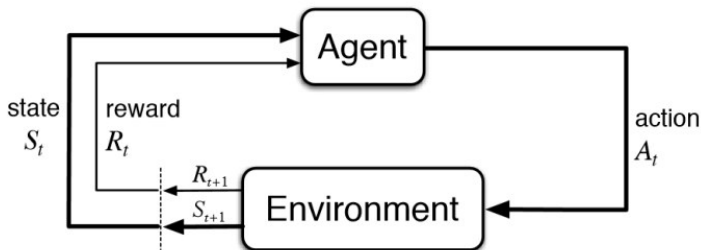
$$V(s) = \underbrace{R(s)}_{\text{Immediate reward}} + \underbrace{\gamma \sum_{s' \in S} P(s'|s) V(s')}_{\text{Discounted sum of future reward}}$$

Iterative Algorithm for Computing Value of a MRP

Algorithm 1 Iterative algorithm to calculate MRP value function

```
1: for all states  $s \in S$ ,  $V'(s) \leftarrow 0$ ,  $V(s) \leftarrow \infty$ 
2: while  $\|V - V'\| > \epsilon$  do
3:    $V \leftarrow V'$ 
4:   For all states  $s \in S$ ,  $V'(s) = R(s) + \gamma \sum_{s' \in S} P(s'|s)V(s')$ 
5: end while
6: return  $V'(s)$  for all  $s \in S$ 
```

Markov Decision Process (MDP)



- 1 Markov Decision Process describes the framework of reinforcement learning. MDP is a tuple: (S, A, P, R, γ)

- 1 S is a finite set of states, A is a finite set of actions
- 2 P^a is dynamics/transition model for each action
 $P(s_{t+1} = s' | s_t = s, a_t = a)$
- 3 R is a reward function $R(s_t = s, a_t = a) = \mathbb{E}[r_t | s_t = s, a_t = a]$
- 4 Discount factor $\gamma \in [0, 1]$

Return Function and Value Function

1 Definition of Horizon

- 1 Number of maximum time steps in each episode
- 2 Can be infinite, otherwise called finite Markov (reward) Process

2 Definition of Return

- 1 Discounted sum of rewards from time step t to horizon

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots + \gamma^{T-t-1} R_T$$

3 Definition of state value function $V_t(s)$ for a MRP

- 1 Expected return from t in state s

$$\begin{aligned} V_t(s) &= \mathbb{E}[G_t | s_t = s] \\ &= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T | s_t = s] \end{aligned}$$

- 2 Present value of future rewards

Value function for MDP

- ① The state-value function $v_\pi(s)$ of an MDP is the expected return starting from state s , and following policy π

$$v_\pi(s) = \mathbb{E}_\pi[G_t | s_t = s] \quad (1)$$

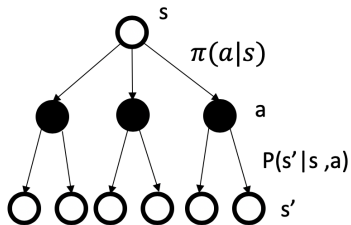
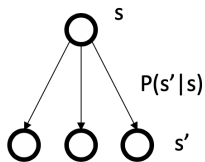
- ② The action-value function $q_\pi(s, a)$ is the expected return starting from state s , taking action a , and then following policy π

$$Q_\pi(s, a) = \mathbb{E}_\pi[G_t | s_t = s, A_t = a] \quad (2)$$

- ③ We have the relation between $v_\pi(s)$ and $Q_\pi(s, a)$

$$v_\pi(s) = \sum_{a \in A} \pi(a|s) Q_\pi(s, a) \quad (3)$$

Comparison of Markov Process and MDP



Bellman Expectation Equation

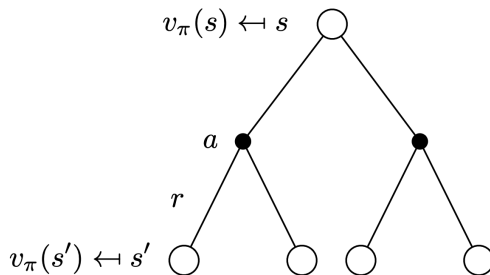
- 1 The state-value function can be decomposed into immediate reward plus discounted value of the successor state,

$$V_{\pi}(s) = E_{\pi}[R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s] \quad (4)$$

- 2 The action-value function can similarly be decomposed

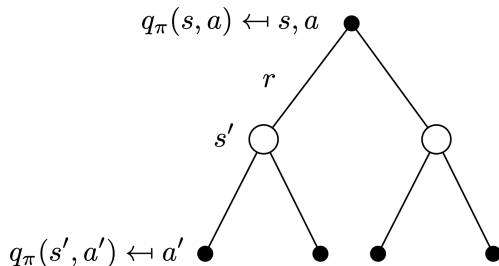
$$Q_{\pi}(s, a) = E_{\pi}[R_{t+1} + \gamma Q_{\pi}(s_{t+1}, A_{t+1}) | s_t = s, A_t = a] \quad (5)$$

Backup Diagram for V^π



$$v_\pi(s) = \sum_{a \in A} \pi(a|s) (R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) v_\pi(s')) \quad (6)$$

Backup Diagram for Q^π



$$q_\pi(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) \sum_{a' \in A} \pi(a'|s') q_\pi(s', a') \quad (7)$$

Decision Making in MDP

- ① Prediction: evaluate a given policy
 - ① Input: MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ and policy π
 - ② Output: value function v_π
- ② Control: search the optimal policy
 - ① Input: MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$
 - ② Output: optimal value function v_* and optimal policy π_*
- ③ Prediction and control can be solved by dynamic programming.

Policy evaluation: Iteration on Bellman expectation backup

Bellman expectation backup for a particular policy

Bootstrapping: estimate on top of estimates

$$v_{k+1}(s) = \sum_{a \in \mathcal{A}} \pi(a|s) (R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) v_k(s')) \quad (8)$$

1 Synchronous backup algorithm:

1 At each iteration $k+1$

update $v_{k+1}(s)$ from $v_k(s')$ for all states $s \in \mathcal{S}$ where s' is a successor state of s

$$v_{k+1}(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) (R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) v_k(s')) \quad (9)$$

2 Convergence: $v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_\pi$

MDP Control

- 1 Compute the optimal policy

$$\pi_*(s) = \arg \max_{\pi} v_{\pi}(s) \quad (10)$$

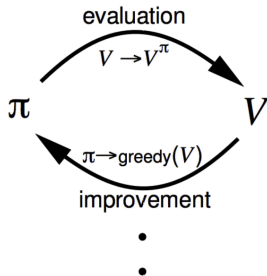
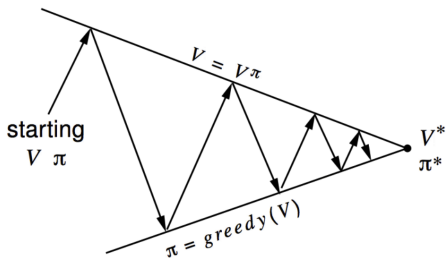
- 2 Policy iteration and value iteration

Improve a Policy through Policy Iteration

① Iterate through the two steps:

- ① **Evaluate** the policy π (computing v given current π)
- ② **Improve** the policy by acting greedily with respect to v_π

$$\pi' = \text{greedy}(v_\pi) \quad (11)$$



Policy Improvement

- 1 compute the state-action value of a policy π :

$$q_{\pi_i}(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) v_{\pi_i}(s') \quad (12)$$

- 2 Compute a new greedy policy π_{i+1} for all $s \in \mathcal{S}$ following

$$\pi_{i+1}(s) = \arg \max_a q_{\pi_i}(s, a) \quad (13)$$

- 3 We can prove that Monotonic Improvement in Policy.

Monotonic Improvement and Bellman optimality equation

- 1 If improvements stop,

$$q_{\pi}(s, \pi'(s)) = \max_{a \in \mathcal{A}} q_{\pi}(s, a) = q_{\pi}(s, \pi(s)) = v_{\pi}(s)$$

- 2 Thus the Bellman optimality equation has been satisfied

$$v_{\pi}(s) = \max_{a \in \mathcal{A}} q_{\pi}(s, a)$$

- 3 Then we have the following Bellman optimality backup:

$$v_{*}(s) = \max_a (R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) v_{*}(s'))$$

Value iteration by turning the Bellman optimality equation as update rule

- ① If we know the solution to subproblem $v_*(s')$, which is optimal.
- ② Then the solution for the optimal $v_*(s)$ can be found by iteration over the following Bellman Optimality backup rule,

$$v(s) \leftarrow \max_{a \in \mathcal{A}} \left(R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) v(s') \right)$$

- ③ The idea of value iteration is to apply these updates iteratively

Summary for Prediction and Control in MDP

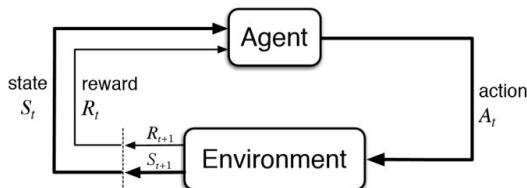
Table: Dynamic Programming Algorithms

Problem	Bellman Equation	Algorithm
Prediction	Bellman Expectation Equation	Iterative Policy Evaluation
Control	Bellman Expectation Equation	Policy Iteration
Control	Bellman Optimality Equation	Value Iteration

Model-free RL

- ① Policy evaluation, policy iteration and value iteration for solving a **known** MDP
- ② Model-free prediction: Estimate value function of an **unknown** MDP
- ③ Model-free control: Optimize value function of an **unknown** MDP

Model-free RL: prediction and optimal control by interaction



- 1 No more oracle or shortcut of the known transition dynamics and reward function
- 2 Trajectories/episodes are collected by the interaction of the agent and the environment
- 3 Each trajectory/episode contains $\{S_1, A_1, R_1, S_2, A_2, R_2, \dots, S_T, A_T, R_T\}$

Model-free prediction: policy evaluation without the access to the model

- ① Estimating the expected return of a particular policy if we don't have access to the MDP models
 - ① Monte Carlo policy evaluation
 - ② Temporal Difference (TD) learning

Monte-Carlo Policy Evaluation

- ① To evaluate state $v(s)$
 - ① Every time-step t that state s is visited in an episode,
 - ② Increment counter $N(s) \leftarrow N(s) + 1$
 - ③ Increment total return $S(s) \leftarrow S(s) + G_t$
 - ④ Value is estimated by mean return $v(s) = S(s)/N(s)$
- ② By law of large numbers, $v(s) \rightarrow v_\pi(s)$ as $N(s) \rightarrow \infty$

Incremental MC Updates

- 1 Collect one episode $(S_1, A_1, R_1, \dots, S_t)$
- 2 For each state s_t with computed return G_t

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

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

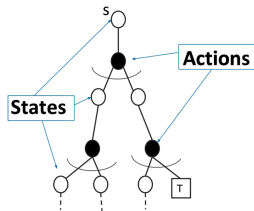
- 3 Or use a running mean (old episodes are forgotten). Good for non-stationary problems.

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

Difference between DP and MC for policy evaluation

- 1 DP computes v_i by bootstrapping the rest of the expected return by the value estimate v_{i-1}
- 2 Iteration on Bellman expectation backup:

$$v_i(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) \left(R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) v_{i-1}(s') \right)$$



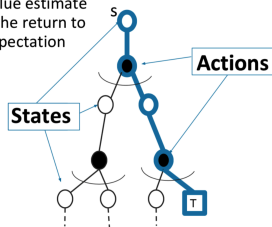
 = Expectation
 = **Terminal state**

Difference between DP and MC for policy evaluation

- 1 MC updates the empirical mean return with one sampled episode

$$v(S_t) \leftarrow v(S_t) + \alpha(G_{i,t} - v(S_t))$$

MC updates the value estimate using a **sample** of the return to approximate an expectation



 = Expectation
 = **Terminal state**

Advantages of MC over DP

- ① MC works when the environment is unknown
- ② Working with sample episodes has a huge advantage, even when one has complete knowledge of the environment's dynamics, for example, transition probability is complex to compute
- ③ Cost of estimating a single state's value is independent of the total number of states. So you can sample episodes starting from the states of interest then average returns

Temporal-Difference (TD) Learning

- ① TD methods learn directly from episodes of experience
- ② TD is model-free: no knowledge of MDP transitions/rewards
- ③ TD learns from **incomplete** episodes, by bootstrapping

Temporal-Difference (TD) Learning

- ① Objective: learn v_π online from experience under policy π
- ② Simplest TD algorithm: TD(0)
 - ① Update $v(S_t)$ toward estimated return $R_{t+1} + \gamma v(S_{t+1})$

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

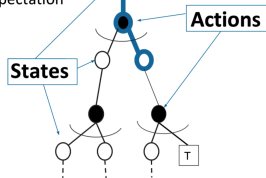
- ③ $R_{t+1} + \gamma v(S_{t+1})$ is called TD target
- ④ $\delta_t = R_{t+1} + \gamma v(S_{t+1}) - v(S_t)$ is called the TD error
- ⑤ Comparison: Incremental Monte-Carlo
 - ① Update $v(S_t)$ toward actual return G_t given an episode i

$$v(S_t) \leftarrow v(S_t) + \alpha(G_{i,t} - v(S_t))$$

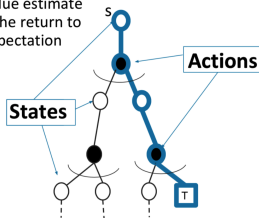
Advantages of TD over MC

TD updates the value estimate using a **sample** of s_{t+1} to approximate an expectation

TD updates the value estimate by **bootstrapping**, uses estimate of $V(s_{t+1})$



MC updates the value estimate using a **sample** of the return to approximate an expectation

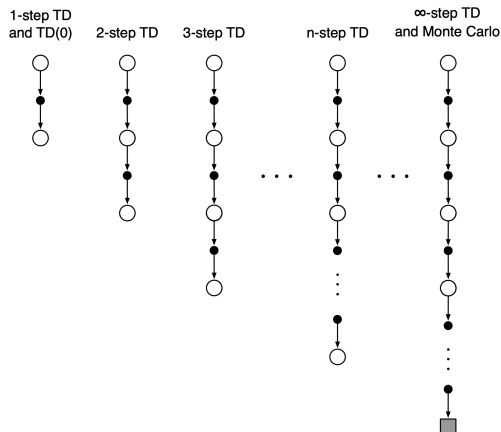


⌋ = Expectation

⌊T⌋ = Terminal state

n-step TD

- 1 n -step TD methods that generalize both one-step TD and MC.
- 2 We can shift from one to the other smoothly as needed to meet the demands of a particular task.



n-step TD prediction

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

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

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

$$\vdots$$

$$n = \infty(MC) \quad G_t^\infty = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T$$

- 2 Thus the n -step return is defined as

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

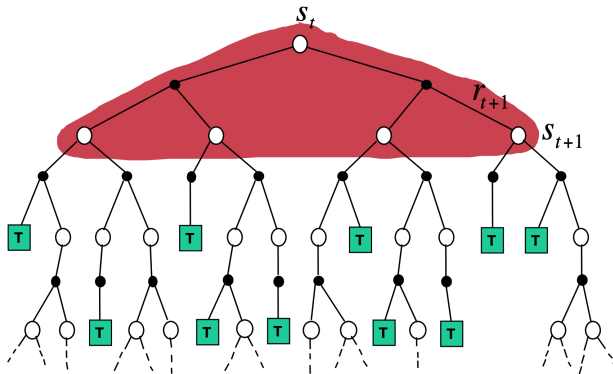
- 3 n -step TD: $v(S_t) \leftarrow v(S_t) + \alpha (G_t^n - v(S_t))$

Bootstrapping and Sampling for DP, MC, and TD

- ① Bootstrapping: update involves an estimate
 - ① MC does not bootstrap
 - ② DP bootstraps
 - ③ TD bootstraps
- ② Sampling: update samples an expectation
 - ① MC samples
 - ② DP does not sample
 - ③ TD samples

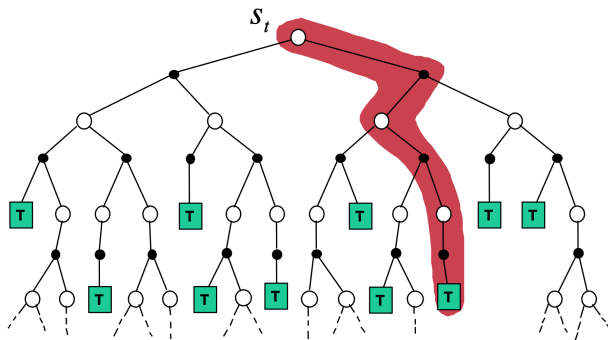
Unified View: Dynamic Programming Backup

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



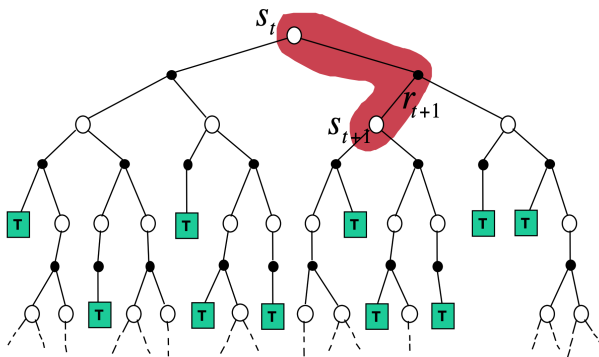
Unified View: Monte-Carlo Backup

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

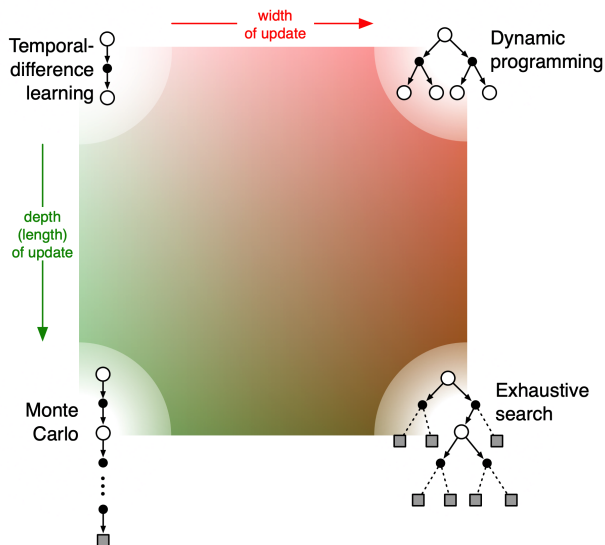


Unified View: Temporal-Difference Backup

$$TD(0) : v(S_t) \leftarrow v(S_t) + \alpha(R_{t+1} + \gamma v(s_{t+1}) - v(S_t))$$



Unified View of Reinforcement Learning



Control: Optimize the policy for a MDP

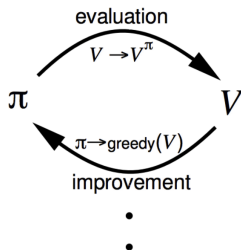
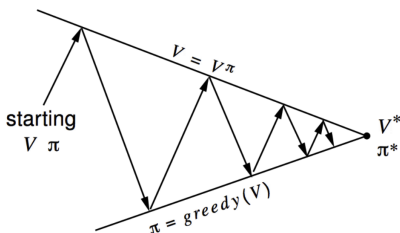
- ① Model-based control: optimize the value function with known MDP
- ② Model-free control: optimize the value function with unknown MDP
- ③ Many model-free RL examples: Go, robot locomotion, patient treatment, helicopter control, Atari, Starcraft

Policy Iteration

1 Iterate through the two steps:

- 1 **Evaluate** the policy π (computing v given current π)
- 2 **Improve** the policy by acting greedily with respect to v_π

$$\pi' = \text{greedy}(v_\pi) \quad (14)$$



Policy Iteration for a Known MDP

- 1 compute the state-action value of a policy π :

$$q_{\pi_i}(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) v_{\pi_i}(s')$$

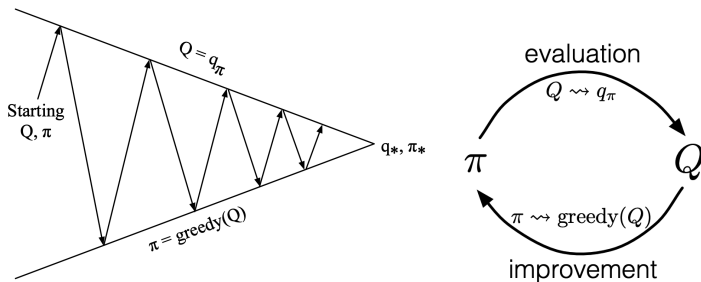
- 2 Compute new policy π_{i+1} for all $s \in \mathcal{S}$ following

$$\pi_{i+1}(s) = \arg \max_a q_{\pi_i}(s, a) \quad (15)$$

- 3 Problem: What to do if there is neither $R(s, a)$ nor $P(s'|s, a)$ known/available?

Generalized Policy Iteration

Monte Carlo version of policy iteration



- 1 Policy evaluation: Monte-Carlo policy evaluation $Q = q_\pi$
- 2 Policy improvement: Greedy policy improvement?

$$\pi(s) = \arg \max_a q(s, a)$$

ϵ -Greedy Exploration

- ❶ Trade-off between exploration and exploitation (we will talk about this in later lecture)
- ❷ ϵ -Greedy Exploration: Ensuring continual exploration
 - ❶ All actions are tried with non-zero probability
 - ❷ With probability $1 - \epsilon$ choose the greedy action
 - ❸ With probability ϵ choose an action at random

$$\pi(a|s) = \begin{cases} \epsilon/|\mathcal{A}| + 1 - \epsilon & \text{if } a^* = \arg \max_{a \in \mathcal{A}} Q(s, a) \\ \epsilon/|\mathcal{A}| & \text{otherwise} \end{cases}$$

Monte Carlo with ϵ -Greedy Exploration

Algorithm 2

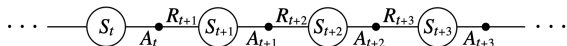
```
1: Initialize  $Q(S, A) = 0, N(S, A) = 0, \epsilon = 1, k = 1$ 
2:  $\pi_k = \epsilon$ -greedy( $Q$ )
3: loop
4:   Sample  $k$ -th episode  $(S_1, A_1, R_2, \dots, S_T) \sim \pi_k$ 
5:   for each state  $S_t$  and action  $A_t$  in the episode do
6:      $N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$ 
7:      $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)}(G_t - Q(S_t, A_t))$ 
8:   end for
9:    $k \leftarrow k + 1, \epsilon \leftarrow 1/k$ 
10:   $\pi_k = \epsilon$ -greedy( $Q$ )
11: end loop
```

TD for Prediction and Control

- ① We can also use TD instead of MC in our control loop
 - ① Apply TD to $Q(S, A)$
 - ② Use ϵ -greedy policy improvement
 - ③ Update every time-step rather than at the end of one episode

Recall: TD Prediction

- 1 An episode consists of an alternating sequence of states and state-action pairs:



- 2 TD(0) method for estimating the value function $V(S)$

$A_t \leftarrow$ action given by π for S

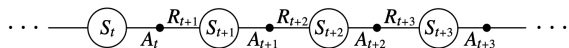
Take action A_t , observe R_{t+1} and S_{t+1}

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

- 3 How about estimating action value function $Q(S)$?

Sarsa: On-Policy TD Control

- 1 An episode consists of an alternating sequence of states and state-action pairs:



- 2 ϵ -greedy policy for one step, then bootstrap the action value function:

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

- 3 The update is done after every transition from a nonterminal state S_t
- 4 TD target $\delta_t = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$

Sarsa algorithm

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

Initialize $Q(s, a)$, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

 Initialize S

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

 Repeat (for each step of episode):

 Take action A , observe R, S'

 Choose A' from S' using policy derived from Q (e.g., ϵ -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

n -step Sarsa

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

$$n = 1(\text{Sarsa}) q_t^{(1)} = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$$

$$n = 2 \quad q_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 Q(S_{t+2}, A_{t+2})$$

\vdots

$$n = \infty(\text{MC}) \quad q_t^\infty = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T$$

- 2 Thus the n -step Q-return is defined as

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

- 3 n -step Sarsa updates $Q(s,a)$ towards the n -step Q-return:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left(q_t^{(n)} - Q(S_t, A_t) \right)$$

Tomorrow

- ① Lecture on on-policy learning and off-policy learning