

# Week 2: Markov Decision Process

Bolei Zhou

*UCLA*

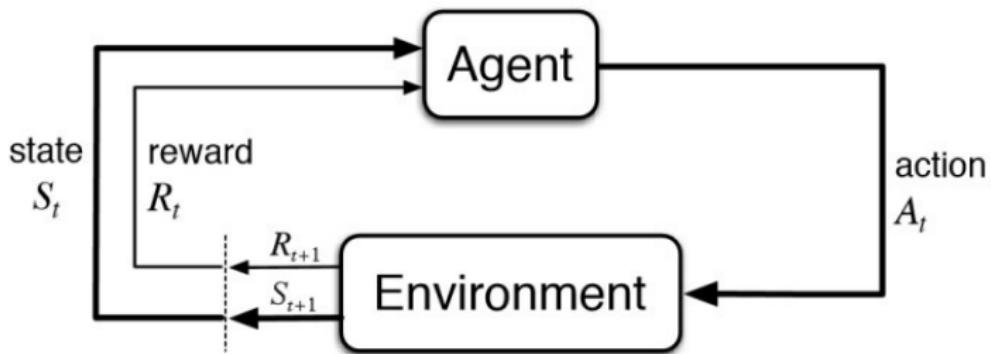
January 13, 2026

# Announcement

- ① Our Google Cloud Education Credits are approved.
  - ① Each enrolled student will give \$50 credit to use.
  - ② Instruction will be sent out next week
  - ③ TAs will give a tutorial in one of the discussions later
- ② Assignment 1 will be due by the end of next week

- ① Last Week
  - ① Course overview
- ② Basic components of RL: reward, policy, value function, action-value function, model
- ③ A simplified RL task: *k*-armed Bandit Problem
- ④ Markov Decision Process (MDP)
  - ① Markov Chain → Markov Reward Process (MRP) → Markov Decision Processes (MDP)
  - ② Policy evaluation in MDP
- ⑤ Control in MDP: policy iteration and value iteration
  - ① Improving dynamic programming
- ⑥ Textbook of Sutton and Barto: Chapter 3 and Chapter 4

# Reinforcement Learning: sequential decision making



Reward: part of environment (pre-defined function), defining reward will change how the agent behaves  
reward hacking

# Reward

- ① A reward is a scalar feedback signal
- ② Reward indicates how well agent is doing at step  $t$
- ③ The objective of RL is to maximize the cumulative reward
  - ① Cumulative reward:  $G_t = \sum_k^{\infty} R_{t+k+1}$        $= 1/N \sum_n G_n$
  - ② Expected cumulative reward:  $E_{\pi}[G_t]$       given this policy, how much can we expect
- ④ Example: roll a dice five times and receive the face number as the reward, what is the expected cumulative reward at the beginning? or before the third time?

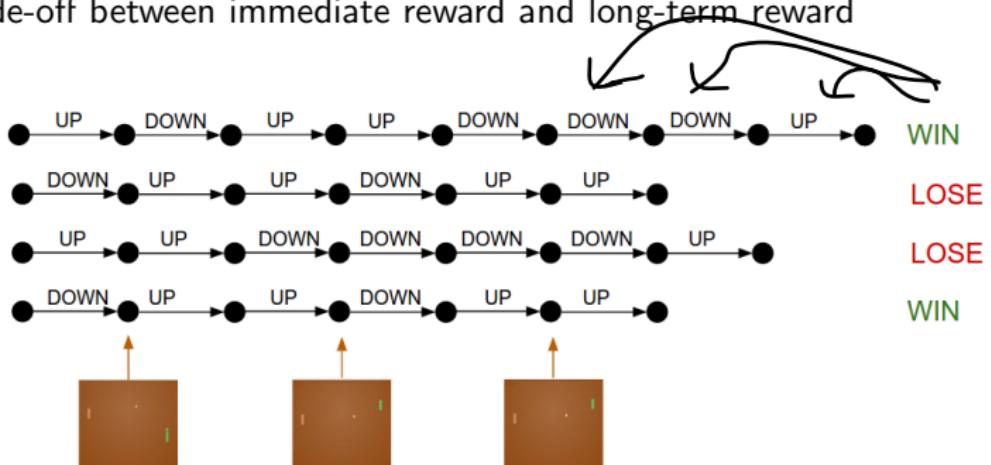
# Sequential Decision Making

- ① Objective of the agent: select a series of actions to maximize the cumulative reward
- ② Actions may have long-term consequences

given a win, we want to reward the actions taken to get to that step

- ① Reward may be delayed

- ② Trade-off between immediate reward and long-term reward



# Major components of an RL agent

- ① An RL agent may include one or more of these components:
  - ① Policy  $\pi(s)$ : agent's behavior function
  - ② Value function  $V(s)$ : how good is each state or action
  - ③ Action-value function  $Q(s,a)$ : how good is an action in a certain state
  - ④ Model: agent's state representation of the environment

# Policy

for classification, policy might be the only thing we care about

same environment + state, we can different actions (b/c distribution)

- ① A policy is the agent's behavior model
- ② It is a map function from state/observation to action
- ③ Stochastic policy: action probability  $\pi(a|s) = P[A_t = a|S_t = s]$
- ④ Deterministic policy: action  $a_t^* = \arg \max_a Q(s_t, a)$

softmax policy

# Value function

- ① Value function: expected discounted sum of future rewards under a particular policy  $\pi$ 
  - ① It is used to quantify the goodness of states if following a particular policy  $\pi$

if we play multiple times (b/c stochastic), we will just take the average  
 $1/N \sum_n G_n$  (otherwise, for deterministic, we don't really need an average)

gamma is an exponential discount (< 1) factor  
- if 1, then future rewards are weighed the same as cur rewards

$$v_\pi(s) = E_\pi[G_t | S_t = s] = E_\pi\left[\sum_k \gamma^k R_{t+k+1} | S_t = s\right] \quad (1)$$

$k \rightarrow \infty$  (until end of episode)

- ② Action-value function, or Q function

$$Q(s, a) = E_\pi[G_t | S_t = s, A_t = a] = E_\pi\left[\sum_k \gamma^k R_{t+k+1} | S_t = s, A_t = a\right]$$

action x state table

- we want to fill in this table with the rewards
- this is the goal of q-learning

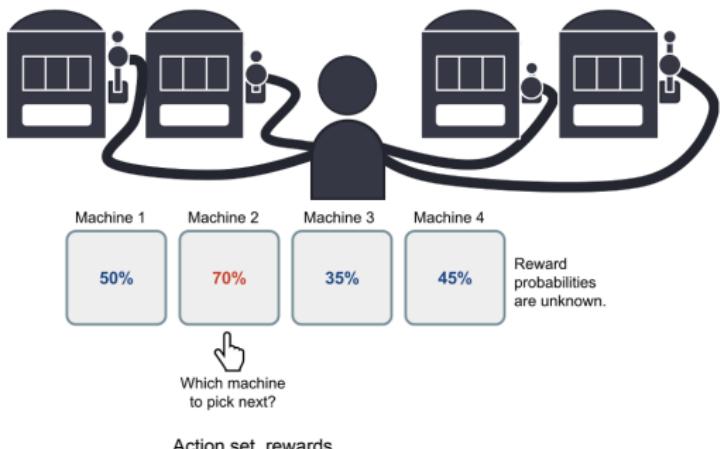
(2)

# Model

A model predicts what the environment will do next  
- think physics

- ① A model predicts what the environment will do next
  - ① Predict the next state:  $P[S_{t+1} = s' | S_t = s, A_t = a]$
  - ② Predict the next reward:  $E[R_{t+1} | S_t = s, A_t = a]$
- ② Sometimes the world model is given, like Newtonian force  $F = ma$

# k-Armed Bandit Problem as a simplified RL example



- ① A multi-armed bandit is a tuple  $\langle \mathcal{A}, \mathcal{R} \rangle$
- ②  $k$  actions to take at each step  $t$
- ③  $\mathcal{R}^a(r) = P(r|a)$  is unknown probability distribution over rewards
- ④ At each step  $t$  the agent selects an action  $a_t \in \mathcal{A}$ , then the environment generates a reward  $r_t \sim \mathcal{R}^{a_t}$
- ⑤ The goal of agent is to maximize cumulative reward  $\sum_{\tau=1}^T r_\tau$

```
class BernoulliArm():    E[G] = P
    def __init__(self, p):
        self.p = p
    def draw(self):
        if random.random() > self.p:
            return 0.0
        else:
            return 1.0

class NormalArm():          E[G] = mu
    def __init__(self, mu, sigma):
        self.mu = mu
        self.sigma = sigma
    def draw(self):
        return random.gauss(self.mu, self.sigma)
```

# Definition of Value Function and Action-Value Function

since there is only one state, the q table is just a vector for the actions reward

- ① The action-value is the mean reward for action  $a$

$$Q(a) = \mathbb{E}(r|a) \quad (3)$$

- ② The optimal value

$$V^* = Q(a^*) = \max_{a \in \mathcal{A}} Q(a) \quad (4)$$

to estimate  $Q(a)$  based on what we have seen prior

- ③ To estimate  $Q(a)$ , we can compute  $Q_t(a)$  at step  $t$  as

$$Q_t(a) = \frac{\text{sum of rewards when } a \text{ taken prior to } t}{\text{number of times } a \text{ taken prior to } t} = \frac{\sum_{i=1}^{t-1} r_i \cdot 1_{A_i=a}}{\sum_{i=1}^{t-1} 1_{A_i=a}} \quad (5)$$

# Greedy Action and $\epsilon$ -Greedy Action to Take

- ① The estimation of  $Q(a)$  at step  $t$

$$Q_t(a) = \frac{\text{sum of rewards when } a \text{ taken prior to } t}{\text{number of times } a \text{ taken prior to } t} = \frac{\sum_{i=1}^{t-1} r_i \cdot 1_{A_i=a}}{\sum_{i=1}^{t-1} 1_{A_i=a}} \quad (6)$$

greedy action: we've learned the q function already, let's exploit and get the arg max

- ② Greedy action selection algorithm:  $A_t = \arg \max_a Q_t(a)$
- ③ Problem with the greedy algorithm?
- ④  $\epsilon$ -Greedy: greedy most of the time, but with small probability  $\epsilon$  select random actions ( $\epsilon$  is usually as 0.1)
  - ① probability  $1 - \epsilon$ :  $A_t = \arg \max_a Q_t(a)$  inference time we use greedy
  - ② probability  $\epsilon$ :  $A_t = \text{uniform}(\mathcal{A})$  gradually reduce probability for exploration from 1 to 0

# $\epsilon$ -Greedy Algorithm

---

## Algorithm 1 A simple *epsilon*-Greedy bandit algorithm

---

```
1: for  $a = 1$  to  $k$  do
2:    $Q(a) = 0$ ,  $N(a) = 0$ 
3: end for
4: loop
5:    $A = \begin{cases} \arg \max_a Q(a) & \text{with probability } 1 - \epsilon \\ \text{uniform}(\mathcal{A}) & \text{with probability } \epsilon \end{cases}$ 
6:    $r = \text{bandit}(A)$  collect reward
7:    $N(A) = N(A) + 1$  update our count for that action
8:    $Q(A) = Q(A) + \frac{1}{N(A)}[r - Q(A)]$ 
9: end loop update our reward function with the residual
```

---

Deriving: NewEstimate = OldEstimate + StepSize[Target - OldEstimate]

$$Q_t(a_t) = Q_{t-1} + \frac{1}{N_t(a_t)}(r_t - Q_{t-1}(a_t)) \quad (7)$$

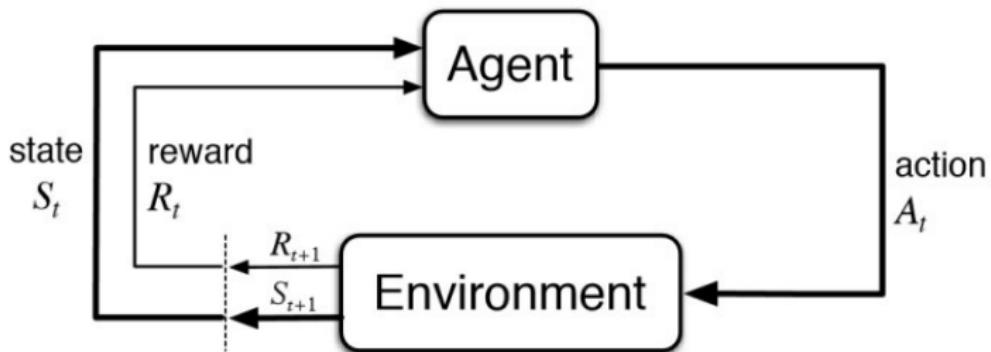
# Why it is a simplified RL task

action leads directly to reward

- ① no delayed reward: instance feedback
- ② no change of state: only one state
- ③ independent of consecutive behaviors

# Markov Decision Process (MDP)

"formal definition"



- ① Markov Decision Process can model a lot of real-world problems. It formally describes the framework of reinforcement learning
- ② Under MDP, the environment is fully observable.
  - ① Optimal control primarily deals with continuous MDPs
  - ② Partially observable problems can be converted into MDPs

# Defining Three Markov Models

- Markov Processes
- Markov Reward Processes (MRPs)
- Markov Decision Processes (MDPs)

# Markov Property

current state has all the information about how the world will move into the next state

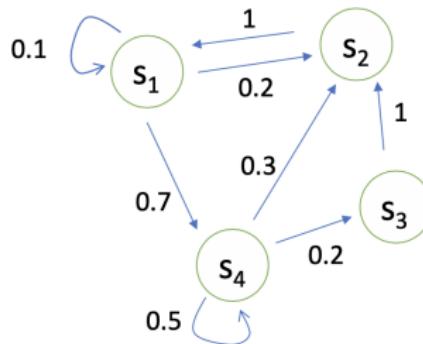
- ① The history of states:  $h_t = \{s_1, s_2, s_3, \dots, s_t\}$
- ② State  $s_t$  is Markovian if and only if:

$$p(s_{t+1}|s_t) = p(s_{t+1}|h_t) \quad (8)$$

$$p(s_{t+1}|s_t, a_t) = p(s_{t+1}|h_t, a_t) \quad (9)$$

- ③ “The future is independent of the past given the present”

# Markov Process/Markov Chain

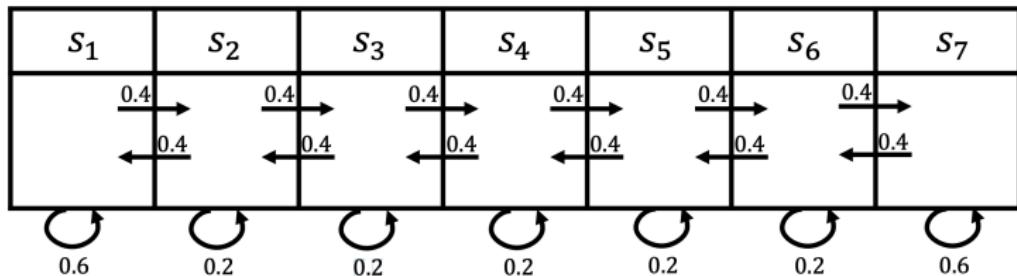


so at  $s_1 \rightarrow$  theres 0.1 we stay, 0.2 we go to  $s_2$ , 0.7 we go to  $s_4$

- ① State transition matrix  $P$  specifies  $p(s_{t+1} = s' | s_t = s)$

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

# Example of MP



① Sample episodes starting from  $s_3$

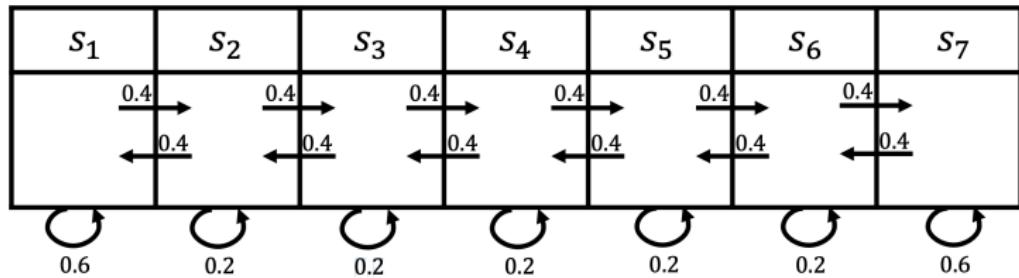
- ①  $s_3, s_4, s_5, s_6, s_6$
- ②  $s_3, s_2, s_3, s_2, s_1$
- ③  $s_3, s_4, s_4, s_5, s_5$

# 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**  $R(s_t = s) = \mathbb{E}[r_t | s_t = s]$  vector of size  $|S|$  (number of states)
  - ④ Discount factor  $\gamma \in [0, 1]$
- ③ If finite number of states,  $R$  can be a vector
  - ① We focus on tabular representation first
  - ② What happen if there are infinite number of states?

for now, assume we have a discrete number of states

# Example of MRP



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 and Value Function

typically, need some shorter term rewards for learning

## ① Definition of Horizon

- ① Number of maximum time steps in each episode/trajectory
- ② Can be infinite, otherwise called finite Markov (reward) Process
- ③ Per game: 100 moves for Go, 80 moves for chess

## 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  $s$

if we sample many times, we can set the state value function  
as  $V(s_{-4}) = 1/n \sum G(s_{-4})$

$$V_t(s) = \mathbb{E}[G_t | s_t = s]$$

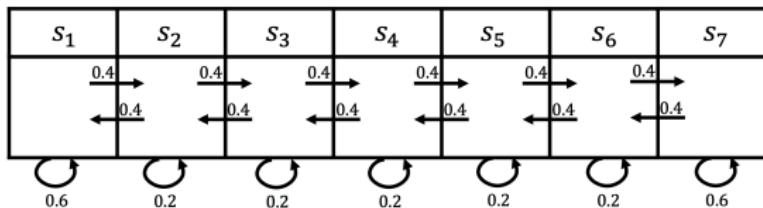
$$= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T | s_t = s]$$

- ② Present value of accumulated future rewards

# Why Discount Factor $\gamma$

- ① Avoid infinite returns in cyclic Markov processes
- ② Uncertainty about the future
  - ① If the reward is financial, immediate rewards may earn more interest than delayed rewards
- ③ Animal/human behavior shows a preference for immediate reward
- ④ It is sometimes possible to use undiscounted Markov reward processes (i.e.  $\gamma = 1$ ), e.g. if all sequences terminate.
  - ①  $\gamma = 0$ : Only care about the immediate reward
  - ②  $\gamma = 1$ : Future reward is equal to the immediate reward.

# Example of MRP



- ① Reward:  $+5$  in  $s_1$ ,  $+10$  in  $s_7$ ,  $0$  in all other states. So that we can represent  $R = [5, 0, 0, 0, 0, 0, 10]$
- ② Sample returns  $G$  for a 3-step episodes with  $\gamma = 1/2$ 
  - ① return for  $s_4, s_5, s_6, s_7$  :  $0 + \frac{1}{2} \times 0 + \frac{1}{4} \times 10 = 2.5$
  - ② return for  $s_4, s_3, s_2, s_1$  :  $0 + \frac{1}{2} \times 0 + \frac{1}{4} \times 5 = 1.25$
  - ③ return for  $s_4, s_5, s_6, s_6$ :  $0$
- ③ How to compute the value function? For example, the value of state  $s_4$  as  $V(s_4) = \mathbb{E}[G_t | s_t = s_4]$

# 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**:

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

- ③ Practice: To derive the Bellman equation from the definition of  $V(s)$

① Hint:  $V(s) = \mathbb{E}[R_{t+1} + \gamma \mathbb{E}[R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots] | s_t = s]$

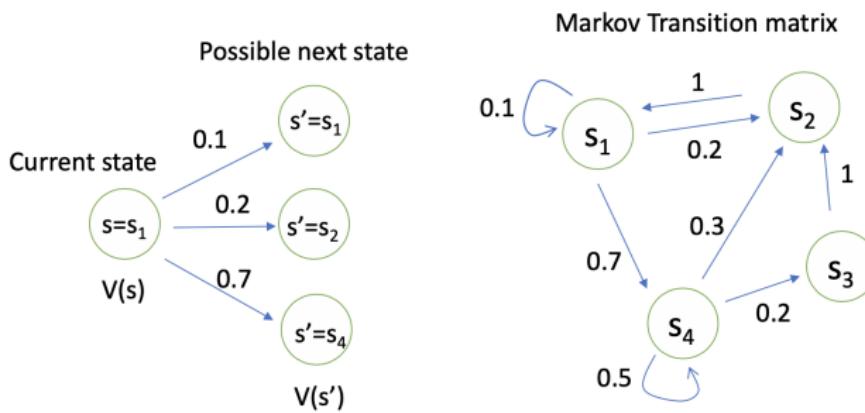
Reward from next

$V(S')$

# Understanding Bellman Equation

- ① **Bellman equation** describes the iterative relations of states

$$V(s) = R(s) + \gamma \sum_{s' \in S} P(s'|s) V(s')$$



# Matrix Form of Bellman Equation for MRP

Therefore, we can express  $V(s)$  using the matrix form:

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

$$V = R + \gamma PV$$

- ① Analytic solution for value of MRP:  $V = (I - \gamma P)^{-1}R$ 
  - ① Matrix inverse takes the complexity  $O(N^3)$  for N states
  - ② Only possible for small MRPs
- ② Other methods to solve this?

# Iterative Algorithm for Computing Value of a MRP

- ① Dynamic Programming      approx
- ② Monte-Carlo evaluation      sampling
- ③ Temporal-Difference learning

between dynamic programming and monte carlo sampling

# Monte Carlo Algorithm for Computing Value of a MRP

---

## Algorithm 2 Monte Carlo simulation to calculate MRP value function

---

```
1:  $i \leftarrow 0, G_t \leftarrow 0$  just sampling with the defined transitions  
collect the returns, and use the discounting, then update using the average  
2: while  $i \neq N$  do  
3:   generate an episode, starting from state  $s$  and time  $t$   
4:   Using the generated episode, calculate return  $g = \sum_{i=t}^{H-1} \gamma^{i-t} r_i$   
5:    $G_t \leftarrow G_t + g, i \leftarrow i + 1$   
6: end while  
7:  $V_t(s) \leftarrow G_t/N$ 
```

---

- ① For example: to calculate  $V(s_4)$  we can generate a lot of trajectories then take the average of the returns:
  - ① return for  $s_4, s_5, s_6, s_7 : 0 + \frac{1}{2} \times 0 + \frac{1}{4} \times 10 = 2.5$
  - ② return for  $s_4, s_3, s_2, s_1 : 0 + \frac{1}{2} \times 0 + \frac{1}{4} \times 5 = 1.25$
  - ③ return for  $s_4, s_5, s_6, s_6 : 0$
  - ④ more trajectories

# Iterative Algorithm for Computing Value of a MRP

Two copies for updates  
 $V(S) =$   
 $V'(S) =$

if there's a large difference for the updates ( $> \text{eps}$ ), then we keep updating  
we then use  $V$  (prior estimate) to update  $V'$  (next estimate)

---

## Algorithm 3 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:      $V \leftarrow V'$   $R(s)$  is a defined reward function
  - 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  $s \in S$
- 

$V$  = value function

this is kind of like "until convergence"

# Markov Decision Process (MDP)

includes an agent with actions

- ① Markov Decision Process is Markov Reward Process with decisions.
- ② Definition of MDP
  - ①  $S$  is a finite set of states
  - ②  **$A$  is a finite set of actions**
  - ③  $P$  is dynamics/transition model for each action  $a$  under certain state  $s$  as  $P(s_{t+1} = s' | s_t = s, a_t = a)$  condition on action too
  - ④  $R$  is a reward function  $R(s_t = s, a_t = a) = \mathbb{E}[r_t | s_t = s, a_t = a]$
  - ⑤ Discount factor  $\gamma \in [0, 1]$  Reward is a 2d table (state and action)
- ③ MDP is a tuple:  $(S, A, P, R, \gamma)$

# Policy in MDP

- ① Policy specifies what action to take in each state
- ② Given a state, specify a distribution over actions
- ③ Policy:  $\pi(a|s) = P(a_t = a|s_t = s)$
- ④ Policies are stationary (time-independent),  $A_t \sim \pi(a|s)$  for any  $t > 0$

input state -> output probability for action

# Policy in MDP

- ① Given a MDP  $(S, A, P, R, \gamma)$  and a policy  $\pi$
- ② The state and reward sequence  $S_1, R_2, S_2, R_2, \dots$  is a Markov reward process  $(S, P^\pi, R^\pi, \gamma)$  where,

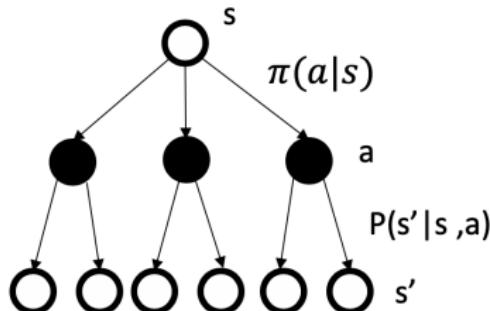
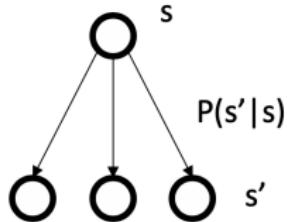
$$P^\pi(s'|s) = \sum_{a \in A} \pi(a|s) P(s'|s, a)$$

$$R^\pi(s) = \sum_{a \in A} \pi(a|s) R(s, a)$$

marginalization based on state and action

# Comparison of MP/MRP and MDP

no agent so no actions



transitions depend on action  
action taken depends on policy  
next state is stochastic and depends on action and state

# 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  $\pi$

$$v^\pi(s) = \mathbb{E}_\pi[G_t | s_t = s] \quad (10)$$

- ② The action-value function  $q^\pi(s, a)$  is the expected return starting from state  $s$ , taking action  $a$ , and then following policy  $\pi$   
**want agent to learn a good q**

$$q^\pi(s, a) = \mathbb{E}_\pi[G_t | s_t = s, A_t = a] \quad (11)$$

- ③ We have the relation between  $v^\pi(s)$  and  $q^\pi(s, a)$

$$v^\pi(s) = \sum_{a \in A} \pi(a|s) q^\pi(s, a) \quad (12)$$

value function is like the expected value at a state  $s$   
using q function as the expectation for  $(s, a)$

# Bellman Expectation Equation

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

$$v^\pi(s) = E_\pi[R_{t+1} + \gamma v^\pi(s_{t+1}) | s_t = s] \quad (13)$$

- ② The action-value function can similarly be decomposed

$$q^\pi(s, a) = E_\pi[R_{t+1} + \gamma q^\pi(s_{t+1}, A_{t+1}) | s_t = s, A_t = a] \quad (14)$$

# Bellman Expectation Equation for $V^\pi$ and $Q^\pi$

$$v^\pi(s) = \sum_{a \in A} \pi(a|s) q^\pi(s, a) \quad (15)$$

$$q^\pi(s, a) = R_s^a + \gamma \sum_{s' \in S} P(s'|s, a) v^\pi(s') \quad (16)$$

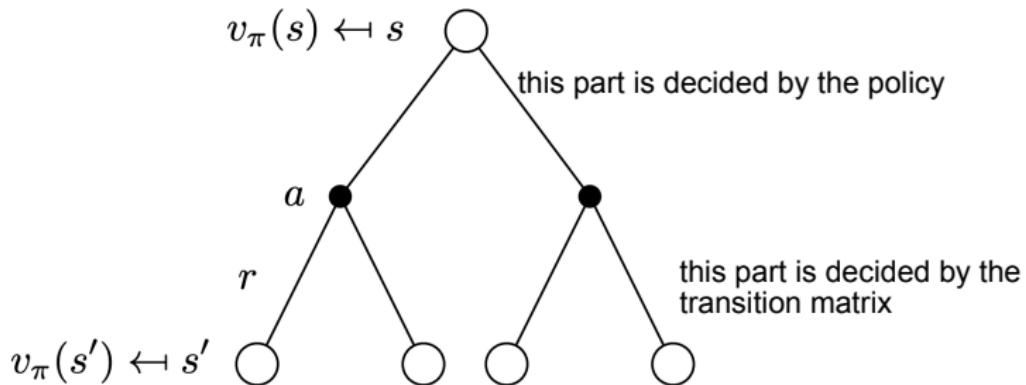
Thus how we marginalize future states

$$v^\pi(s) = \sum_{a \in A} \pi(a|s) (R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) v^\pi(s')) \quad (17)$$

how we marginalize an action that occurs

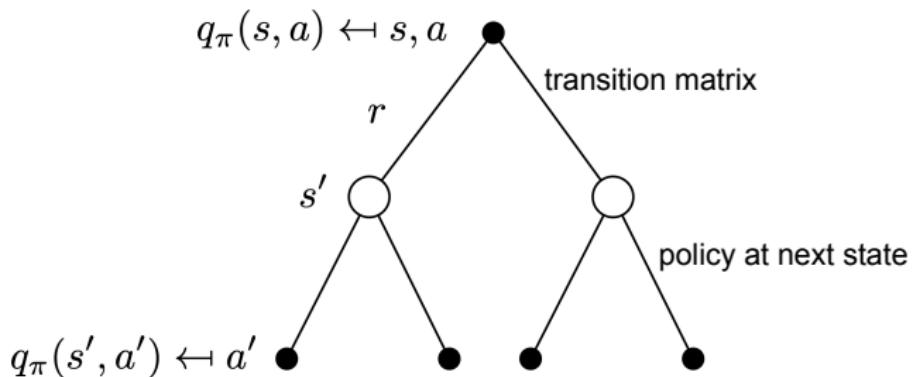
$$q^\pi(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) \sum_{a' \in A} \pi(a'|s') q^\pi(s', a') \quad (18)$$

# Backup Diagram for $V^\pi$



$$v^\pi(s) = \sum_{a \in A} \pi(a|s)(R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a)v^\pi(s')) \quad (19)$$

# Backup Diagram for $Q^\pi$



$$q^\pi(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) \sum_{a' \in A} \pi(a'|s') q^\pi(s', a') \quad (20)$$

# Policy Evaluation

state of value given a fixed policy

- ① Evaluate the value of state given a policy  $\pi$ : compute  $v^\pi(s)$
- ② Also called as (value) prediction

## Example: Navigate the boat



Figure: Markov Chain/MRP: Go with river stream



Figure: MDP: Navigate the boat

## Example: Policy Evaluation

$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$	$s_7$
						

- ① Two actions: *Left* or *Right*
- ② For all actions, reward: +5 in  $s_1$ , +10 in  $s_7$ , 0 in all other states. So that we can represent  $R = [5, 0, 0, 0, 0, 0, 10]$
- ③ Let's have a deterministic policy  $\pi(s) = \text{Left}$  and  $\gamma = 0$  for any state  $s$ , then what is the value of the policy?
  - ①  $V^\pi = [5, 0, 0, 0, 0, 0, 10]$  since  $\gamma = 0$

# Example: Policy Evaluation

$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$	$s_7$
				$V_0 = [5, 0, 0, 0, 0, 0, 10]$	$V_1 = [5, 0.5 \times 5, 0, 0, 0, 0, 10]$	$V_2 = [5, 2.5, 0]$

Learning      January 1

- ①  $R = [5, 0, 0, 0, 0, 0, 10]$
- ② Practice 1: Deterministic policy  $\pi(s) = Left$  with  $\gamma = 0.5$  for any state  $s$ , then what are the state values under the policy?
- ③ Practice 2: Stochastic policy  $P(\pi(s) = Left) = 0.5$  and  $P(\pi(s) = Right) = 0.5$  and  $\gamma = 0.5$  for any state  $s$ , then what are the state values under the policy?
- ④ Iteration t:  
 $v_t^\pi(s) = \sum_a P(\pi(s) = a)(r(s, a) + \gamma \sum_{s' \in S} P(s'|s, a)v_{t-1}^\pi(s'))$

# Decision Making in Markov Decision Process (MDP)

- ① Prediction (evaluate a given policy):
  - ① Input: MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$  and policy  $\pi$  or MRP  $\langle \mathcal{S}, \mathcal{P}^\pi, \mathcal{R}^\pi, \gamma \rangle$
  - ② Output: value function  $v^\pi$
- 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  $\pi^*$
- ③ Prediction and control in MDP can be solved by dynamic programming.

# Dynamic Programming

Dynamic Programming is a very general solution method for problems which have two properties:

- ① Optimal substructure

- ① Principle of optimality applies
  - ② Optimal solution can be decomposed into subproblems

- ② Overlapping subproblems

- ① Subproblems recur many times
  - ② Solutions can be cached and reused

Markov decision processes satisfy both properties

- ① Bellman equation gives recursive decomposition
- ② Value function stores and reuses solutions

# Prediction: Policy evaluation on MDP

- ① Objective: Evaluate a given policy  $\pi$  for a MDP
- ② Output: the value function under policy  $v^\pi$
- ③ Solution: iteration on Bellman expectation backup
- ④ Algorithm: Synchronous backup
  - ① At each iteration  $t+1$   
update  $v_{t+1}(s)$  from  $v_t(s')$  for all states  $s \in \mathcal{S}$  where  $s'$  is a successor state of  $s$   
where do we get transitions? defined by env (sometimes hidden from agent)

$$v_{t+1}(s) = \sum_{a \in \mathcal{A}} \pi(a|s)(R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \hat{v}_t(s')) \quad (21)$$

previous estimate

- ⑤ Convergence:  $v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v^\pi$

## Policy evaluation: Iteration on Bellman expectation backup

Bellman expectation backup for a particular policy

$$v_{t+1}(s) = \sum_{a \in \mathcal{A}} \pi(a|s)(R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a)v_t(s')) \quad (22)$$

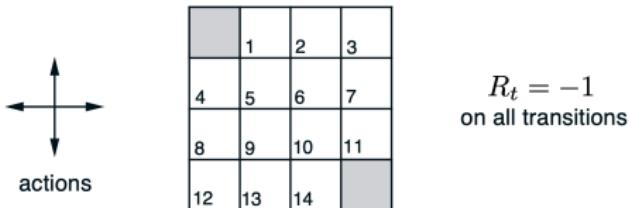
Or if in the form of MRP  $\langle \mathcal{S}, \mathcal{P}^\pi, \mathcal{R}, \gamma \rangle$

marginalize actions?

$$v_{t+1}(s) = R^\pi(s) + \gamma \sum_{s' \in \mathcal{S}} P^\pi(s'|s)v_t(s') \quad (23)$$

# Evaluating a Random Policy in the Small Gridworld

Example 4.1 in the Sutton RL textbook.



- ① Undiscounted episodic MDP ( $\gamma = 1$ )
- ② Nonterminal states 1, ..., 14
- ③ Two terminal states (two shaded squares)
- ④ Action leading out of grid leaves state unchanged,  $P(7|7, \text{right}) = 1$
- ⑤ Reward is  $-1$  until the terminal state is reached
- ⑥ Transition is deterministic given the action, e.g.,  $P(6|5, \text{right}) = 1$
- ⑦ Uniform random policy  $\pi(l|.) = \pi(r|.) = \pi(u|.) = \pi(d|.) = 0.25$

# Evaluating a Random Policy in the Small Gridworld

- ① Iteratively evaluate the random policy

$v_k$  for the  
Random Policy

$k = 0$

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

$k = 3$

0.0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0.0

-1 is defined by the reward function

-1 reward except end state

$k = 1$

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

$k = 10$

0.0	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	0.0

$k = 2$

0.0	-1.7	-2.0	-2.0
-1.7	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.7
-2.0	-2.0	-1.7	0.0

$k = \infty$

0.0	-14.	-20.	-22.
-14.	-18.	-20.	-20.
-20.	-20.	-18.	-14.
-22.	-20.	-14.	0.0

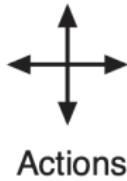
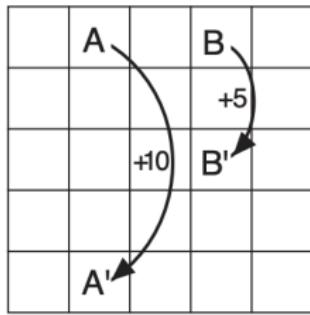
# A live demo on policy evaluation

$$v^\pi(s) = \sum_{a \in A} \pi(a|s)(R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a)v^\pi(s')) \quad (24)$$

- ① [https://cs.stanford.edu/people/karpathy/reinforcejs/  
gridworld\\_dp.html](https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld_dp.html)

# Practice: Gridworld

Textbook Example 3.5: GridWorld



3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

# MDP Control

- ① Compute the optimal policy

$$\pi^*(s) = \arg \max_{\pi} v^{\pi}(s) \quad (25)$$

- ② Optimal policy for a MDP in an infinite horizon problem (agent acts forever) is

- ① Deterministic
- ② Stationary (does not depend on time step)
- ③ Unique? Not necessarily, may have state-actions with identical optimal values

# Optimal Value Function

- ① The optimal state-value function  $v^*(s)$  is the maximum value function over all policies

$$v^*(s) = \max_{\pi} v^{\pi}(s)$$

- ② The optimal policy

$$\pi^*(s) = \arg \max_{\pi} v^{\pi}(s)$$

- ③ An MDP is “solved” when we know the optimal value
- ④ There exists a unique optimal value function, but could be multiple optimal policies (two actions that have the same optimal value function)

# Finding Optimal Policy

if we have an optimal q function, then we can just take the argmax action at each state

- ① An optimal policy can be found by maximizing over  $q^*(s, a)$ ,

$$\pi^*(a|s) = \begin{cases} 1, & \text{if } a = \arg \max_{a \in A} q^*(s, a) \\ 0, & \text{otherwise} \end{cases}$$

- ② There is always a deterministic optimal policy for any MDP
- ③ If we know  $q^*(s, a)$ , we immediately have the optimal policy

# Policy Search

using the brute force algo to find the best policy

- ① One option is to enumerate search the best policy
- ② Number of deterministic policies is  $|\mathcal{A}|^{\mathcal{|S|}}$  can take  $|\mathcal{A}|$  actions at each state
- ③ Other approaches such as policy iteration and value iteration are more efficient
  - ① Policy iteration
  - ② Value iteration

# Improving a Policy through Policy Iteration

- under current policy, we evaluate the policy
- then we act greedy

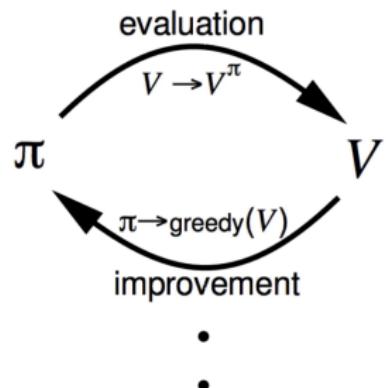
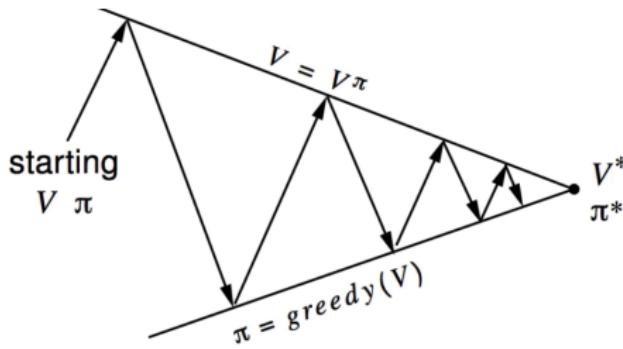
① Iterate through the two steps: keep doing this and update the q table

① Evaluate the policy  $\pi$  (computing  $v$  given current  $\pi$ )

② Improve the policy by acting greedily with respect to  $v^\pi$

$$\boxed{\begin{aligned} v^\pi(s) &= \sum_{a \in A} \pi(a|s) q^\pi(s, a) \\ q^\pi(s, a) &= R_s^a + \gamma \sum_{s' \in S} P(s'|s, a) v^\pi(s') \end{aligned}}$$

$$\pi' = \text{greedy}(v^\pi) \quad (26)$$



we have the relationship between  $v$  and  $q$  with the bellman expectation equations

# Policy Improvement

- ① Compute the state-action value of a policy  $\pi$ :

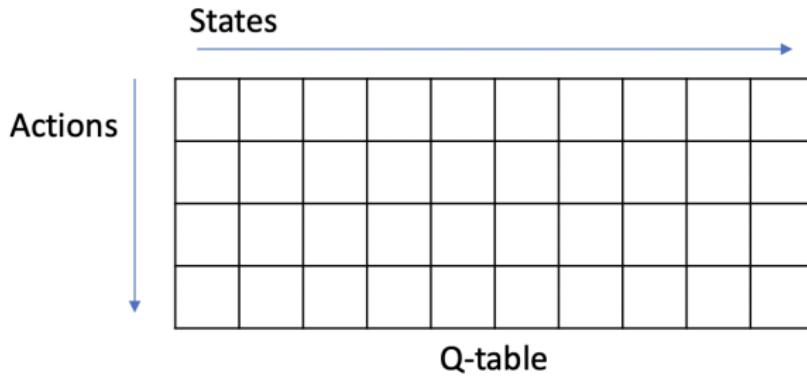
$$q^{\pi_i}(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) v^{\pi_i}(s') \quad (27)$$

until convergence given the current policy

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

$$\pi_{i+1}(s) = \arg \max_a q^{\pi_i}(s, a) \quad (28)$$

then we update



# Monotonic Improvement in Policy

- ① Consider a deterministic policy  $a = \pi(s)$
- ② We improve the policy through

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

- ③ This improves the value from any state  $s$  over one step,

$$q^\pi(s, \pi'(s)) = \max_{a \in \mathcal{A}} q^\pi(s, a) \geq q^\pi(s, \pi(s)) = v^\pi(s)$$

- ④ It therefore improves the value function,  $v^{\pi'}(s) \geq v^\pi(s)$   
updated policy will always be better

$$\begin{aligned} v^\pi(s) &\leq q^\pi(s, \pi'(s)) = \mathbb{E}_{\pi'}[R_{t+1} + \gamma v^\pi(S_{t+1}|S_t = s)] \\ &\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma q^\pi(S_{t+1}, \pi'(S_{t+1}))|S_t = s] \\ &\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 q^\pi(S_{t+2}, \pi'(S_{t+2}))|S_t = s] \\ &\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \dots | S_t = s] = v^{\pi'}(s) \end{aligned}$$

# Monotonic Improvement in Policy

since we have monotonic improvement, we have the best possible policy once the q function stops improving

- 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 Therefore  $v^\pi(s) = v^*(s)$  for all  $s \in \mathcal{S}$ , so  $\pi$  is an optimal policy

# Bellman Optimality Equation

- ① The optimal value functions are reached by the Bellman optimality equations:

this is not argmax, we just take the largest  $q(s, a)$

$$v^*(s) = \max_a q^*(s, a)$$

$$q^*(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a)v^*(s')$$

thus (plugging in  $v^*(s)$  into  $q^*(s, a)$  equation

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

$$q^*(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) \max_{a'} q^*(s', a')$$

## 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

# Algorithm of Value Iteration

- ① Objective: find the optimal policy  $\pi$
- ② Solution: iteration on the Bellman optimality backup
- ③ Value Iteration algorithm:

- ① initialize  $k = 1$  and  $v_0(s) = 0$  for all states  $s$
- ② For  $k = 1 : H$

① for each state  $s$  just need to keep a set of values

update q

$$q_{k+1}(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a)v_k(s') \quad (29)$$

update v

$$v_{k+1}(s) = \max_a q_{k+1}(s, a) \quad (30)$$

②  $k \leftarrow k + 1$

- ③ To retrieve the optimal policy after the value iteration:

$$\pi(s) = \arg \max_a R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a)v_{k+1}(s') \quad (31)$$

Type text here

# Example: Shortest Path

g			

Problem

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$v_1$

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1

$v_2$

0	-1	-2	-2
-1	-2	-2	-2
-2	-2	-2	-2
-2	-2	-2	-2

$v_3$

0	-1	-2	-3
-1	-2	-3	-3
-2	-3	-3	-3
-3	-3	-3	-3

$v_4$

0	-1	-2	-3
-1	-2	-3	-4
-2	-3	-4	-4
-3	-4	-4	-4

$v_5$

0	-1	-2	-3
-1	-2	-3	-4
-2	-3	-4	-5
-3	-4	-5	-5

$v_6$

0	-1	-2	-3
-1	-2	-3	-4
-2	-3	-4	-5
-3	-4	-5	-6

$v_7$

After the optimal values are reached, we run policy extraction to retrieve the optimal policy.

# Difference between Policy Iteration and Value Iteration

- ① Policy iteration includes: **policy evaluation + policy improvement**, and the two are repeated iteratively until policy converges.
- ② Value iteration includes: **finding optimal value function + one policy extraction**. There is no repeat of the two because once the value function is optimal, then the policy out of it should also be optimal (i.e. converged).
- ③ Finding optimal value function can also be seen as a combination of policy improvement (due to max) and truncated policy evaluation (the reassignment of  $v(s)$  after just one sweep of all states regardless of convergence).

# 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

# Demo of policy iteration and value iteration

0.22 ↑	0.25 ↑	0.27 ↑	0.30 ↑	0.34 ↑	0.38 ↓	0.34 ↑	0.30 ↓	0.34 ↑	0.38 ↓
0.25 →	0.27 →	0.30 →	0.34 →	0.38 →	0.42 ↓	0.38 →	0.34 ↔	0.38 →	0.42 ↓
0.21 ↓					0.46 ↓				0.46 ↓
0.20 ↗	0.22 ↑	0.25 ↓	0.25 ↗ $R_{-1.0}$		0.52 →	0.57 →	0.64 ↓	0.57 ↗	0.52 ↗
0.22 ↑	0.25 ↑	0.27 ↓	0.25 ↗		0.08 ↓ $R_{-1.0}$	-0.36 → $R_{-1.0}$	0.71 ↓	0.64 →	0.57 →
0.25 ↑	0.27 ↑	0.30 ↓	0.27 ↗		1.20 ↑ $R_{-1.0}$	0.08 → $R_{-1.0}$	0.79 ↓	-0.29 → $R_{-1.0}$	0.52 ↓
0.27 ↑	0.30 ↑	0.34 ↓	0.30 ↔		1.0 ↑ $R_{-1.0}$	0.97 → $R_{-1.0}$	0.87 ↓	-0.21 → $R_{-1.0}$	0.57 ↓
0.31 ↑	0.34 ↑	0.38 ↓	0.38 ↗ $R_{-1.0}$		-0.03 ↓ $R_{-1.0}$	-0.03 ↑ $R_{-1.0}$	0.71 ↑	0.71 →	0.64 →
0.34 →	0.38 →	0.42 →	0.46 →	0.52 →	0.57 →	0.84 →	0.71 ↑	0.64 ↑	0.57 ↑
0.31 ↓	0.34 ↓	0.38 ↓	0.42 ↓	0.46 ↓	0.52 ↓	0.57 ↓	0.64 ↑	0.57 ↑	0.52 ↑

- ① Policy iteration: Iteration of policy evaluation and policy improvement(update)
- ② Value iteration
- ③ [https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld\\_dp.html](https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld_dp.html)

# Policy iteration and value iteration on FrozenLake

① <https://github.com/ucla-rlcourse/RLexample/tree/master/MDP>

# Improving Dynamic Programming

- ① A major drawback to the DP methods is that they involve operations over the entire state set of the MDP, that is, they require sweeps of the state set.
- ② If the state set is very large, for example, the game of backgammon has over  $10^{20}$  states. Thousands of years to be taken to finish one sweep.
- ③ Asynchronous DP algorithms are in-place iterative DP that are not organized in terms of systematic sweeps of the state set
- ④ The values of some states may be updated several times before the values of others are updated once.

# Improving Dynamic Programming

v\_new v\_old, copy values of v\_old into v\_new

Synchronous dynamic programming is usually slow. Three simple ideas to extend DP for asynchronous dynamic programming:

- ① In-place dynamic programming
- ② Prioritized sweeping
- ③ Real-time dynamic programming

# In-Places Dynamic Programming

- ① Synchronous value iteration stores two copies of value function:

for all  $s$  in  $\mathcal{S}$

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

$$v_{old} \leftarrow v_{new}$$

- ② In-place value iteration only stores one copy of value function:

for all  $s$  in  $\mathcal{S}$

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

have only one value vector and update in place

- issue is that some values may be already updated

# Prioritized Sweeping

- ① Use magnitude of Bellman error to guide state selection, e.g.

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

- ② Backup the state with the largest remaining Bellman error
- ③ Update Bellman error of affected states after each backup
- ④ Can be implemented efficiently by maintaining a priority queue

sort so values with larger error are updated first

# Real-Time Dynamic Programming

"don't have to wait for one episode to update"

- ① To solve a given MDP, we can run an iterative DP algorithm at the same time that an agent is actually experiencing the MDP
- ② The agent's experience can be used to determine the states to which the DP algorithm applies its updates
- ③ We can apply updates to states as the agent visits them. So focus on the parts of the state set that are most relevant to the agent
- ④ After each time-step  $S_t, A_t$ , backup the state  $S_t$ ,

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

when agent just moves one step, we update

# Sample Backups

- ① The key design for RL algorithms such as Q-learning and SARSA in next lectures
- ② Using sample rewards and sample transition pairs  $\langle S, A, R, S' \rangle$ , rather than the reward function  $\mathcal{R}$  and transition dynamics  $\mathcal{P}$
- ③ Benefits:
  - ① Model-free: no advance knowledge of MDP required
  - ② Break the curse of dimensionality through sampling
  - ③ Cost of backup is constant, independent of  $n = |\mathcal{S}|$

# Approximate Dynamic Programming

- ① Using a function approximator  $\hat{v}(s, \mathbf{w})$
- ② Fitted value iteration repeats at each iteration  $k$ ,
  - ① Sample state  $s$  from the state cache  $\tilde{\mathcal{S}}$

$$\tilde{v}_k(s) = \max_{a \in \mathcal{A}} \left( R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \hat{v}(s', \mathbf{w}_k) \right)$$

- ② Train next value function  $\hat{v}(s', \mathbf{w}_{k+1})$  using targets  $\langle s, \tilde{v}_k(s) \rangle$ .
- ③ Key idea behind the Deep Q-Learning

- ① Summary: MDP, policy evaluation, policy iteration, and value iteration
- ② Check out the code example of MDP at <https://github.com/ucla-rlcourse/RExample/tree/master/MDP>
- ③ In this week's discussion session, TAs will give you more examples
- ④ Next Week: Model-free methods
- ⑤ Reading: Textbook Chapter 5 and Chapter 6