

# Dynamic Programming

## Reinforcement Learning : Module 04

Policy Evaluation (Prediction), Policy Improvement, Policy Iteration, Value Iteration, Asynchronous Dynamic Programming, Generalized Policy Iteration

The slides contain

- Excerpts from the book Reinforcement Learning : An Introduction, 2<sup>nd</sup> edition, Richard S. Sutton and Andrew G. Barto
- Content acquired using [ChatGPT](#), [Gemini](#) and CoPilot

# Dynamic Programming

- The term dynamic programming (DP) refers to a collection of algorithms that can be used to compute optimal policies given a perfect model of the environment as a Markov decision process (MDP).
- Classical DP algorithms are of limited utility in reinforcement learning both because of their assumption of a perfect model and because of their great computational expense, but they are still important theoretically.
- The methods - to be discussed - can be viewed as attempts to achieve much the same effect as DP, only with less computation and without assuming a perfect model of the environment.

# Dynamic Programming

- The key idea of DP, and of reinforcement learning generally, is the use of value functions to organize and structure the search for good policies.
- We can easily obtain optimal policies once we have found the optimal value functions,  $v_*$  or  $q_*$ , which satisfy the Bellman optimality equations:

$$\begin{aligned} v_*(s) &= \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma v_*(s') \right], \text{ or} \end{aligned} \tag{4.1}$$

$$\begin{aligned} q_*(s, a) &= \mathbb{E} \left[ R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma \max_{a'} q_*(s', a') \right], \end{aligned} \tag{4.2}$$

for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ , and  $s' \in \mathcal{S}^+$ .

# Dynamic Programming

- (As we shall see)
- DP algorithms are obtained by turning Bellman equations such as 4.1 and 4.2 into assignments, that is, into **update rules** for **improving approximations** of the desired value functions.
  - **Update rules** : Imagine these update rules as instructions for progressively improving the agent's estimates of the value functions ( $v$  or  $q$ ).
  - **Approximations and Improvement**: In real-world applications, the environment might be too complex to calculate the exact value functions directly. So, DP algorithms rely on approximations. These approximations are initially rough estimates of the true value functions.

# The Update Process

- The update rules based on the Bellman equations iteratively improve these approximations.
- In each iteration, the agent considers the rewards received, the value of the next state (according to the current approximation), and updates the estimated value of the current state.
- This process continues until the approximations converge to a good estimate of the actual value functions.
- **Analogy** : Imagine passing through a maze blindfolded.

# Policy Evaluation (Prediction)

# Revision : What is Policy?

- Policy defines the strategy an agent uses to navigate its environment and make decisions.
- **Function of a Policy:**
  - The policy acts as a mapping function. It takes the current state of the environment (a set of features representing the situation) as input and outputs an action for the agent to take.
  - Ideally, the policy guides the agent towards actions that maximize its long-term reward.
- **Types of Policies:**
  - **Deterministic Policies:** These policies always recommend the same action for a given state. Imagine a robot following a pre-programmed path in a factory.
  - **Stochastic Policies:** These policies assign a probability distribution over possible actions for each state. The agent randomly chooses an action based on these probabilities. This can be useful for exploration in unknown environments or dealing with uncertainty.

# Policy Evaluation (Prediction)

- Let us consider how to compute the state-value function  $v_\pi$  for an arbitrary policy  $\pi$ .
- This is called policy evaluation in the DP literature.
  - Also referred as the prediction problem.



# Policy Evaluation (Prediction)

- For all  $s \in \mathcal{S}$ ,

$$\begin{aligned} v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \\ &= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \\ &= \mathbb{E}_{\pi}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s] \end{aligned} \tag{4.3}$$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma v_{\pi}(s') \right] \tag{4.4}$$

where

- $G_t$  is expected return, sum of the sequence of rewards received after time step  $t$ .
- $\pi(a|s)$  is the probability of taking action  $a$  in state  $s$  under policy  $\pi$ , and the expectations are subscripted by  $\pi$  to indicate that they are conditional on  $\pi$  being followed

# Iterative Policy Evaluation

- Consider a sequence of approximate value functions  $v_0, v_1, v_2, \dots$ , each mapping  $\mathbf{S}^+$  to  $\mathbf{R}$  (the real numbers).
- The initial approximation,  $v_0$ , is chosen arbitrarily (except that the terminal state, if any, may be given value 0), and each successive approximation is obtained by using the Bellman equation for  $v_\pi$  (4.4) as an update rule:

$$\begin{aligned} v_{k+1}(s) &\doteq \mathbb{E}_\pi[R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_k(s')] \end{aligned} \quad (4.5)$$

- This algorithm is called *iterative policy evaluation*.

# Iterative Policy Evaluation

- To produce each successive approximation,  $v_{k+1}$  from  $v_k$ , iterative policy evaluation applies the same operation to each state  $s$ : it replaces the old value of  $s$  with a new value obtained from the old values of the successor states of  $s$ , and the expected immediate rewards, along all the one-step transitions possible under the policy being evaluated.
- We call this kind of operation an expected update.
- Each iteration of iterative policy evaluation updates the value of every state once to produce the new approximate value function  $v_{k+1}$

# Iterative Policy Evaluation

- There are several different kinds of expected updates, depending on whether a state (as here) or a state–action pair is being updated, and depending on the precise way the estimated values of the successor states are combined.
- All the updates done in DP algorithms are called expected updates because they are based on an expectation over all possible next states rather than on a sample next state.

# Developing a computer program

- To write a sequential computer program to implement iterative policy evaluation as given by (4.5) you would have to use two arrays, one for the old values,  $v_k(s)$ , and one for the new values,  $v_{k+1}(s)$ .
- With two arrays, the new values can be computed one by one from the old values without the old values being changed.
- Of course it is easier to use one array and update the values “in place,” that is, with each new value immediately overwriting the old one.
- Then, depending on the order in which the states are updated, sometimes new values are used instead of old ones on the right-hand side of (4.5).

# Developing a computer program

- This in-place algorithm also converges to  $v_\pi$ ; in fact, it usually converges faster than the two-array version, because it uses new data as soon as they are available.
- We think of the updates as being done in a sweep through the state space.
- For the in-place algorithm, the order in which states have their values updated during the sweep has a significant influence on the rate of convergence.
- We usually have the in-place version in mind when we think of DP algorithms.

# Pseudocode

## Iterative Policy Evaluation, for estimating $V \approx v_\pi$

Input  $\pi$ , the policy to be evaluated

Algorithm parameter: a small threshold  $\theta > 0$  determining accuracy of estimation

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

Loop:

$\Delta \leftarrow 0$

Loop for each  $s \in \mathcal{S}$ :

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until  $\Delta < \theta$

# Pseudocode

- A complete in-place version of iterative policy evaluation is shown in pseudocode in the box.
- Note how it handles termination.
- Formally, iterative policy evaluation converges only in the limit, but in practice it must be halted short of this.
- The pseudocode tests the quantity  $\max_{s \in \mathcal{S}} |v_{k+1}(s) - v_k(s)|$  after each sweep and stops when it is sufficiently small.



# Example 4.1

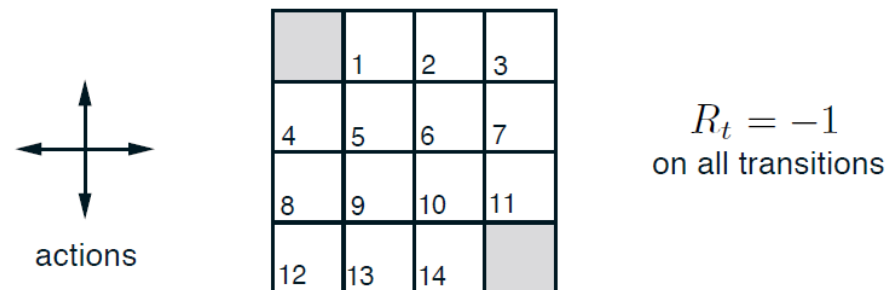
- Consider the 4 X 4 gridworld shown below.



	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

$R_t = -1$   
on all transitions

## Example 4.1



The nonterminal states are  $\mathcal{S} = \{1, 2, \dots, 14\}$ . There are four actions possible in each state,  $\mathcal{A} = \{\text{up}, \text{down}, \text{right}, \text{left}\}$ , which deterministically cause the corresponding state transitions, except that actions that would take the agent off the grid in fact leave the state unchanged. Thus, for instance,  $p(6, -1 | 5, \text{right}) = 1$ ,  $p(7, -1 | 7, \text{right}) = 1$ , and  $p(10, r | 5, \text{right}) = 0$  for all  $r \in \mathcal{R}$ . This is an undiscounted, episodic task. The reward is  $-1$  on all transitions until the terminal state is reached. The terminal state is shaded in the figure (although it is shown in two places, it is formally one state). The expected reward function is thus  $r(s, a, s') = -1$  for all states  $s, s'$  and actions  $a$ . Suppose the agent follows the equiprobable random policy (all actions equally likely). The left side of Figure 4.1 shows the sequence of value functions  $\{v_k\}$  computed by iterative policy evaluation.

# Example 4.1

	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

$R_t = -1$   
on all transitions

$k = 0$

Calculate  $V_1(1)$

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

$v_k$  for the  
random policy

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

greedy policy  
w.r.t.  $v_k$

	↔	↔	↔
↔	↔	↔	↔
↔	↔	↔	↔
↔	↔	↔	

← random  
policy

Action order assumed is : Left,  
Up, Right, Down

$$\begin{aligned} V_1(1) &= [0.25 * (-1 + 1 * v_0(0))] + [0.25 * (-1 + 1 * v_0(1))] + [0.25 * (-1 + 1 * v_0(2))] \\ &\quad + [0.25 * (-1 + 1 * v_0(5))] \\ &= [0.25 * (-1 + 0)] + [0.25 * (-1 + 0)] \\ &\quad + [0.25 * (-1 + 0)] + [0.25 * (-1 + 0)] \\ &= -1 \end{aligned}$$

# Example 4.1

	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

$R_t = -1$   
on all transitions

$k = 0$

$v_k$  for the  
random policy

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

greedy policy  
w.r.t.  $v_k$

	↔	↔	↔
↕	↕	↕	↕
↔	↔	↔	↔
↕	↕	↕	

random  
policy

$k = 1$

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

	←	↔	↔
↑	↕	↕	↕
↔	↔	↔	↓
↕	↕	→	

Calculate  $V_2(1)$

Action order assumed is : Left, Up, Right, Down

$$\begin{aligned}
 V_2(1) &= [0.25 * (-1 + 1 * v_1(0))] + [0.25 * (-1 + 1 * v_1(1))] + [0.25 * (-1 + 1 * v_1(2))] + [0.25 * (-1 + 1 * v_1(5))] \\
 &= [0.25 * (-1 + 0)] + [0.25 * (-1 + -1)] + [0.25 * (-1 + -1)] + [0.25 * (-1 + -1)] \\
 &= -0.25 + (-0.5) + (-0.5) + (-0.5) + (-0.5) \\
 &= -1.75
 \end{aligned}$$

# Example 4.1

$k = 3$

$v_k$  for the  
random policy

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

greedy policy  
w.r.t.  $v_k$

	←	←	↙
↑	↖	↙	↓
↑	↗	↘	↓
↘	→	→	

$k = 10$

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

	←	←	↙
↑	↖	↙	↓
↑	↗	↘	↓
↘	→	→	

optimal  
policy

$k = \infty$

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

	←	←	↙
↑	↖	↙	↓
↑	↗	↘	↓
↘	→	→	

## Example 4.1

- The final estimate is in fact  $v_{\pi}$ , which in this case gives for each state the negation of the expected number of steps from that state until termination.

# Exercise

	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

$R_t = -1$  on all transitions  $k = \infty$

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

- In Example 4.1, if  $\pi$  is the equiprobable random policy,
  - What is  $q_\pi(11, \text{down})$ ?
  - What is  $q_\pi(7, \text{down})$ ?

## Answer

$$q_\pi(s, a) = r + \gamma v_\pi(s')$$

where reward is always -1 and rewards are not discounted.

$$q_\pi(s, a) = -1 + v_\pi(s')$$

$$q_\pi(11, \text{down}) = -1 + v_\pi(15) = -1 + 0 = -1$$

$$q_\pi(7, \text{down}) = -1 + v_\pi(11) = -1 + (-14) = -15 \quad (v_\pi(11) \text{ is looked up from figure 4.1})$$

# Exercise

*Exercise 4.2* In Example 4.1, suppose a new state 15 is added to the gridworld just below state 13, and its actions, **left**, **up**, **right**, and **down**, take the agent to states 12, 13, 14, and 15, respectively. Assume that the transitions *from* the original states are unchanged. What, then, is  $v_\pi(15)$  for the equiprobable random policy? Now suppose the dynamics of state 13 are also changed, such that action **down** from state 13 takes the agent to the new state 15. What is  $v_\pi(15)$  for the equiprobable random policy in this case?  $\square$



# Answer

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

where actions are equiprobable, reward is always -1 and rewards are not discounted.

In first case, the new state is not reachable from state 13.

$$v_{\pi}(s) = \sum_a 0.25[(-1) + v_{\pi}(s')]$$

where  $v(12) = -22, v(13) = -20, v(14) = -14, v(9) = -20$  :

$$v_{\pi}(15) = (-1) + 0.25 * [v_{\pi}(12) + v_{\pi}(13) + v_{\pi}(14) + v_{\pi}(15)]$$

$$v_{\pi}(15) = -15 + 0.25v_{\pi}(15) = -20$$

In second case, the new state is reachable from state 13. Now value of 13 depends on value of 15, and value of 15 depends on value of 13. We have two equations with two unknowns which can be solved.

$v_{\pi}(15)$  is :

$$v_{\pi}(15) = (-1) + 0.25 * [v_{\pi}(12) + v_{\pi}(13) + v_{\pi}(14) + v_{\pi}(15)]$$

$$v_{\pi}(15) = -1 + 0.25 * [-36 + v_{\pi}(13) + v_{\pi}(15)] = -1 - 9 + \frac{v_{\pi}(13)}{4} + \frac{v_{\pi}(15)}{4}$$

$$v_{\pi}(15) = \frac{-40 + v_{\pi}(13)}{3}$$

$v_{\pi}(13)$  is :

$$v_{\pi}(13) = (-1) + 0.25 * [v_{\pi}(12) + v_{\pi}(9) + v_{\pi}(14) + v_{\pi}(15)]$$

$$v_{\pi}(13) = -1 + 0.25 * [-56 + v_{\pi}(15)] = -1 - 14 + \frac{v_{\pi}(15)}{4} = -15 + \frac{v_{\pi}(15)}{4}$$

One can plug  $v_{\pi}(13)$  to obtain  $v_{\pi}(15)$  :

$$v_{\pi}(15) = \frac{-40 + -15 + \frac{v_{\pi}(15)}{4}}{3} = -20$$

# Exercise

- Assume a game to be played as follows.
  1. There are 2 options for the player: *play* or *quit*
  2. If quit is chosen then 10 points are granted to the player and the game ends.
  3. If play is chosen then 4 points are granted to the player and the game continues with the rolling of a 6 faced die.
    - i. If 1 or 2 occurs on the die then the game ends, with no additional points awarded to player.
    - ii. If 3 to 6 occurs on the die then then follow step 1.
- Question : Design this game as MDP.

# Exercise

*Exercise 4.3* What are the equations analogous to (4.3), (4.4), and (4.5) for the action-value function  $q_\pi$  and its successive approximation by a sequence of functions  $q_0, q_1, q_2, \dots$ ?

$$\begin{aligned} q_\pi(s, a) &\doteq \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a] \\ &= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \\ &= \mathbb{E}_\pi\left[R_{t+1} + \gamma \sum_{s', a'} q_\pi(s', a') \mid S_t = s, A_t = a\right] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma \sum_{a'} \pi(a' \mid s') q_\pi(s', a') \right] \end{aligned}$$

$$\begin{aligned} q_{k+1}(s, a) &\doteq \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma \sum_{a'} \pi(a' \mid s') q_k(s', a') \right] \end{aligned}$$

# Policy Improvement

The reason for computing the value function for a policy is to help find better policies.

# What is Policy Improvement

- The process of creating a new policy that outperforms the original policy.
- Achieved by evaluating the current policy and identifying actions that lead to higher rewards in specific states.
- The new policy prioritizes these actions, leading to potentially better outcomes.

# Policy Improvement : Introduction

- Suppose we have determined the value function  $v_\pi$  for an arbitrary deterministic policy  $\pi$ .
- For some state  $s$  we would like to know whether or not we should change the policy to deterministically choose an action  $a \neq \pi(s)$ .
- We know how good it is to follow the current policy from  $s$ —that is  $v_\pi(s)$ —but would it be better or worse to change to the new policy ?
- One way to answer this question is to consider selecting  $a$  in  $s$  and thereafter following the existing policy,  $\pi$ .
- The value of this way of behaving is

$$\begin{aligned} q_\pi(s, a) &\doteq \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma v_\pi(s') \right]. \end{aligned} \tag{4.6}$$

# Policy Improvement : Introduction

$$\begin{aligned} q_{\pi}(s, a) &\doteq \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma v_{\pi}(s') \right]. \end{aligned} \tag{4.6}$$

- The key criterion is whether this is greater than or less than  $v_{\pi}(s)$ .
- If it is greater
  - that is, if it is better to select  $a$  once in  $s$  and thereafter follow  $\pi$  than it would be to follow  $\pi$  all the time
  - then one would expect it to be better still to select  $a$  every time  $s$  is encountered, and that the new policy would in fact be a better one overall.

# Policy Improvement Theorem

Let  $\pi$  and  $\pi'$  be any pair of deterministic policies such that, for all  $s \in \mathcal{S}$ ,

$$q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s). \quad (4.7)$$

Then the policy  $\pi'$  must be as good as, or better than,  $\pi$ . That is, it must obtain greater or equal expected return from all states  $s \in \mathcal{S}$ :

$$v_{\pi'}(s) \geq v_{\pi}(s). \quad (4.8)$$



# Policy Improvement Theorem : Proof

Starting from (4.7), we keep expanding the  $q_\pi$  side with (4.6) and reapplying (4.7) until we get  $v_{\pi'}(s)$ :

$$\begin{aligned} v_\pi(s) &\leq q_\pi(s, \pi'(s)) \\ &= \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = \pi'(s)] && \text{(by (4.6))} \\ &= \mathbb{E}_{\pi'}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s] \\ &\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma q_\pi(S_{t+1}, \pi'(S_{t+1})) \mid S_t = s] && \text{(by (4.7))} \\ &= \mathbb{E}_{\pi'}[R_{t+1} + \gamma \mathbb{E}_{\pi'}[R_{t+2} + \gamma v_\pi(S_{t+2}) \mid S_{t+1}, A_{t+1} = \pi'(S_{t+1})] \mid S_t = s] \\ &= \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 v_\pi(S_{t+2}) \mid S_t = s] \\ &\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 v_\pi(S_{t+3}) \mid S_t = s] \\ &\vdots \\ &\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \cdots \mid S_t = s] \\ &= v_{\pi'}(s). \end{aligned}$$

# Policy Improvement

- We saw how - given a policy and its value function, we can easily evaluate a change in the policy at a single state to a particular action.
- It is a natural extension to consider changes at all states and to all possible actions, selecting at each state the action that appears best according to  $q_\pi(s, a)$ . In other words, to consider the new greedy policy,  $\pi'$ , given by

$$\begin{aligned}\pi'(s) &\doteq \operatorname{argmax}_a q_\pi(s, a) \\ &= \operatorname{argmax}_a \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \operatorname{argmax}_a \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma v_\pi(s') \right],\end{aligned}\tag{4.9}$$

# Policy Improvement

- The greedy policy takes the action that looks best in the short term—after one step of lookahead—according to  $v_\pi$ .
- By construction, the greedy policy meets the conditions of the policy improvement theorem (4.7), so we know that it is as good as, or better than, the original policy.
- The process of making a new policy that improves on an original policy, by making it greedy with respect to the value function of the original policy, is called *policy improvement*.

# Policy Improvement

Suppose the new greedy policy,  $\pi'$ , is as good as, but not better than, the old policy  $\pi$ . Then  $v_\pi = v_{\pi'}$ , and from (4.9) it follows that for all  $s \in \mathcal{S}$ :

$$\begin{aligned} v_{\pi'}(s) &= \max_a \mathbb{E}[R_{t+1} + \gamma v_{\pi'}(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma v_{\pi'}(s') \right]. \end{aligned}$$

But this is the same as the Bellman optimality equation (4.1), and therefore,  $v_{\pi'}$  must be  $v_*$ , and both  $\pi$  and  $\pi'$  must be optimal policies. Policy improvement thus must give us a strictly better policy except when the original policy is already optimal.

# Policy Iteration

# Policy Iteration

Once a policy,  $\pi$ , has been improved using  $v_\pi$  to yield a better policy,  $\pi'$ , we can then compute  $v_{\pi'}$  and improve it again to yield an even better  $\pi''$ . We can thus obtain a sequence of monotonically improving policies and value functions:

$$\pi_0 \xrightarrow{\text{E}} v_{\pi_0} \xrightarrow{\text{I}} \pi_1 \xrightarrow{\text{E}} v_{\pi_1} \xrightarrow{\text{I}} \pi_2 \xrightarrow{\text{E}} \cdots \xrightarrow{\text{I}} \pi_* \xrightarrow{\text{E}} v_*,$$

where  $\xrightarrow{\text{E}}$  denotes a policy *evaluation* and  $\xrightarrow{\text{I}}$  denotes a policy *improvement*. Each policy is guaranteed to be a strict improvement over the previous one (unless it is already optimal). Because a finite MDP has only a finite number of policies, this process must converge to an optimal policy and optimal value function in a finite number of iterations.

This way of finding an optimal policy is called **policy iteration**.

# Policy Iteration Algorithm

**Policy Iteration (using iterative policy evaluation) for estimating  $\pi \approx \pi_*$**

1. Initialization

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

2. Policy Evaluation

Loop:

$\Delta \leftarrow 0$

Loop for each  $s \in \mathcal{S}$ :

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until  $\Delta < \theta$  (a small positive number determining the accuracy of estimation)

3. Policy Improvement

*policy-stable*  $\leftarrow$  *true*

For each  $s \in \mathcal{S}$ :

*old-action*  $\leftarrow \pi(s)$

$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

If *old-action*  $\neq \pi(s)$ , then *policy-stable*  $\leftarrow$  *false*

If *policy-stable*, then stop and return  $V \approx v_*$  and  $\pi \approx \pi_*$ ; else go to 2

## Example 4.2 : Jack's Car Rental

- Students are recommended to study Jack's Car Rental problem discussed on Page 81 of the book.



## Exercise 4.4

- The policy iteration algorithm on page 80 has a subtle bug in that it may never terminate if the policy continually switches between two or more policies that are equally good.
- This is ok for pedagogy, but not for actual use.
- Modify the pseudocode so that convergence is guaranteed.

# Solution

In the step 3. Policy Improvement, it said:

*If old-action  $\neq \pi(s)$  , then .....*

It is a bug and one way to fix it is to say the following instead:

*If old-action  $\notin \{a_i\}$ , which is the all equi-best solutions from  $\pi(s)$ , .....*

## Exercise 4.5

- How would policy iteration be defined for action values?
- Give a complete algorithm for computing  $q_*$ , analogous to that on page 80 for computing  $v_*$ .
- Please pay special attention to this exercise, because the ideas involved will be used throughout the rest of the book.

# Solution

## 1. Initialization

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

## 2. Policy Evaluation

Loop:

$\Delta \leftarrow 0$

Loop for each  $s \in S$  and  $a \in A$ :

$q = Q(s, a)$

$Q(s, a) \leftarrow \sum_{s', r} p(s', r | s, a) \left[ r + \gamma \sum_{a'} \pi(a' | s') Q(s', a') \right]$

$\Delta \leftarrow \max(\Delta, |q - Q(s, a)|)$

until  $\Delta < \theta$  (a small positive number determining the accuracy of estimation)

## 3. Policy Improvement

*policy-stable*  $\leftarrow$  true

For each  $s \in S$  and  $a \in A$ :

*old-action*  $\leftarrow \pi(s)$

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

If *old-action*  $\notin \{a_i\}$ , which is the set of equi-best solutions from  $\pi(s)$

Then *policy-stable*  $\leftarrow$  false

If *policy-stable*, then stop and return  $Q \approx q_*$  and  $\pi \approx \pi_*$ ; else go to 2

# Value Iteration

# Value Iteration

- One drawback to policy iteration is that each of its iterations involves policy evaluation, which may itself be a protracted iterative computation requiring multiple sweeps through the state set.
- If policy evaluation is done iteratively, then convergence exactly to  $v_\pi$  occurs only in the limit.
- Must we wait for exact convergence, or can we stop short of that?
- The example in Figure 4.1 certainly suggests that it may be possible to truncate policy evaluation.
  - In that example, policy evaluation iterations beyond the first three have no effect on the corresponding greedy policy.

# Value Iteration

- In fact, the policy evaluation step of policy iteration can be truncated in several ways without losing the convergence guarantees of policy iteration.
- One important special case is when policy evaluation is stopped after just one sweep (one update of each state).
- This algorithm is called *value iteration*.

# Value Iteration

- It can be written as a particularly simple update operation that combines the policy improvement and truncated policy evaluation steps:

$$\begin{aligned} v_{k+1}(s) &\doteq \max_a \mathbb{E}[R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_k(s')], \end{aligned} \tag{4.10}$$

- Note that value iteration is obtained simply by turning the Bellman optimality equation into an update rule.
- Also note how the value iteration update is identical to the policy evaluation update (4.5) except that it requires the maximum to be taken over all actions



# Value Iteration

- Let us consider how value iteration terminates.
- Like policy evaluation, value iteration formally requires an infinite number of iterations to converge exactly to  $v_*$ .
- In practice, we stop once the value function changes by only a small amount in a sweep.
- The box below shows a complete algorithm with this kind of termination condition.

# Value Iteration

## Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold  $\theta > 0$  determining accuracy of estimation

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

Loop:

```
|  $\Delta \leftarrow 0$   
| Loop for each  $s \in \mathcal{S}$ :  
|    $v \leftarrow V(s)$   
|    $V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$   
|    $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
```

until  $\Delta < \theta$

Output a deterministic policy,  $\pi \approx \pi_*$ , such that

$$\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

## Example 4.3 : Gambler's Problem

- Students are recommended to study Gambler's problem discussed on Page 84 of the book.

# Asynchronous Dynamic Programming

# Asynchronous Dynamic Programming

- A major drawback to the DP methods that we have discussed so far is that they involve operations over the entire state set of the MDP, that is, they require sweeps of the state set.
- If the state set is very large, then even a single sweep can be prohibitively expensive.
- For example, the game of backgammon has over  $10^{20}$  states. Even if we could perform the value iteration update on a million states per second, it would take over a thousand years to complete a single sweep.

# Asynchronous Dynamic Programming

- Asynchronous Dynamic Programming (ADP) refers to a class of algorithms used to solve Markov Decision Processes (MDPs) by asynchronously updating the value function or policy of states.
- Unlike traditional dynamic programming algorithms, such as value iteration and policy iteration, which update all states in a synchronized manner, ADP updates states individually and in any order.
- This asynchronous updating process offers several advantages in reinforcement learning settings:

# Advantages of ADP

- **Efficiency:**

- ADP can be more computationally efficient than synchronous dynamic programming algorithms, especially in large-scale problems where updating all states simultaneously may be computationally prohibitive.
- By focusing computational resources on states that are most in need of updating, ADP can converge to an optimal solution more quickly.

- **Scalability:**

- Asynchronous updates allow for more scalable solutions to reinforcement learning problems, as they enable the algorithm to handle large state spaces more efficiently.
- By updating states independently, ADP can be applied to problems with millions or even billions of states.

# Advantages of ADP

- **Exploration-Exploitation Trade-off:**

- ADP can help balance exploration and exploitation in reinforcement learning. By updating states asynchronously, the algorithm can focus on exploring regions of the state space that are less explored or have higher uncertainty, while exploiting known information in other regions.

- **Incremental Updates:**

- ADP typically performs incremental updates to the value function or policy, updating only the affected states based on changes in neighboring states.
- This incremental update process can lead to faster convergence and reduced computational overhead compared to recomputing the entire value function or policy in each iteration.



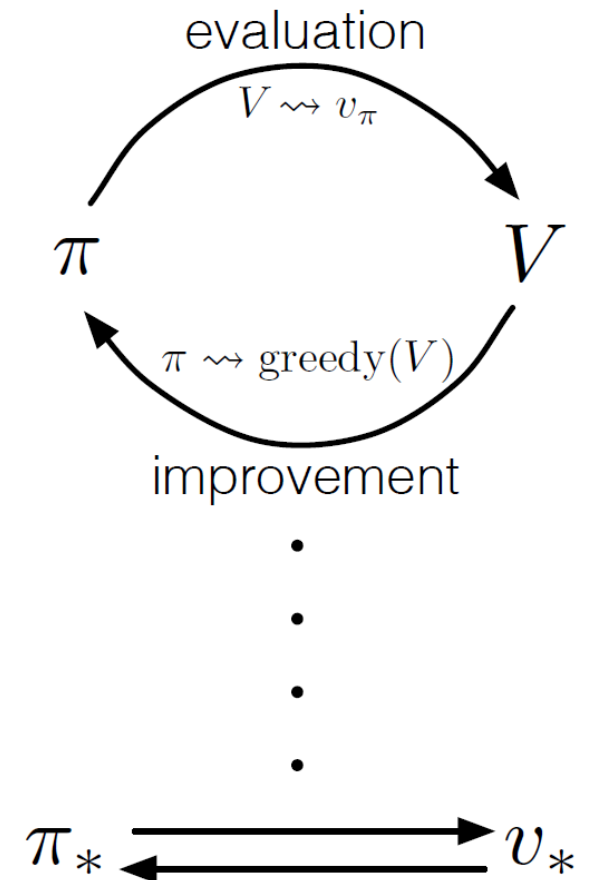
# Asynchronous Dynamic Programming

- Popular ADP algorithms in reinforcement learning include asynchronous value iteration, asynchronous policy iteration, and various variants of asynchronous Q-learning.
- These algorithms have been successfully applied to a wide range of reinforcement learning tasks, including robotic control, game playing, and autonomous decision-making

# Generalized Policy Iteration

# Generalized Policy Iteration

- The term *generalized policy iteration* (GPI) refers to the general idea of letting policy-evaluation and policy improvement processes interact, independent of the granularity and other details of the two processes.
- Almost all reinforcement learning methods are well described as GPI.
- That is, all have identifiable policies and value functions, with the policy always being improved with respect to the value function and the value function always being driven toward the value function for the policy, as suggested by the diagram to the right.



# Generalized Policy Iteration

- If both the evaluation process and the improvement process stabilize, that is, no longer produce changes, then the value function and policy must be optimal.
- The value function stabilizes only when it is consistent with the current policy, and the policy stabilizes only when it is greedy with respect to the current value function.
- Thus, both processes stabilize only when a policy has been found that is greedy with respect to its own evaluation function.
- This implies that the Bellman optimality equation (4.1) holds, and thus that the policy and the value function are optimal.

# Generalized Policy Iteration

- The evaluation and improvement processes in GPI can be viewed as both competing and cooperating.
- They compete in the sense that they pull in opposing directions.
- Making the policy greedy with respect to the value function typically makes the value function incorrect for the changed policy, and making the value function consistent with the policy typically causes that policy no longer to be greedy.
- In the long run, however, these two processes interact to find a single joint solution: the optimal value function and an optimal policy.

# Generalized Policy Iteration

- Generalized Policy Iteration (GPI) is crucial in reinforcement learning due to its ability to iteratively refine both the policy and the value function.
- By continuously evaluating and improving the agent's policy, GPI enables adaptive learning and optimal decision-making in complex environments.
- Its flexibility allows for asynchronous updates, interleaving policy evaluation and improvement, and accommodating various algorithms, making it applicable to a wide range of reinforcement learning problems.
- GPI provides a unified framework that balances exploration and exploitation, leading to faster convergence and more efficient learning.

End