

# NOTE FOR DEEP REINFORCEMENT LEARNING

---

WANG HUI

*Amateur Machine Learning Researcher in NUDT*

[wanghuichn1234@gmail.com](mailto:wanghuichn1234@gmail.com)

---

# Contents

<b>I</b>	<b>Math Foundation</b>	<b>1</b>
<b>1</b>	<b>Basic Concept</b>	<b>1</b>
1.1	Agent-Environment interaction cycle	1
1.2	Episodic task, Continuing task and Return	1
1.3	States	1
1.4	Actions	2
1.5	Transition Function	2
1.6	Rewards	3
1.7	Horizon	3
1.8	Discount Factor	3
<b>II</b>	<b>Balancing immediate and long-term goals</b>	<b>5</b>
<b>2</b>	<b>Policy</b>	<b>5</b>
<b>3</b>	<b>State-value function: What to expect from here?</b>	<b>5</b>
<b>4</b>	<b>Action-value function: What should I expect from here if I do this action?</b>	<b>7</b>
<b>5</b>	<b>Optimal Policy</b>	<b>7</b>
<b>III</b>	<b>Two algorithm to find objects</b>	<b>10</b>
<b>6</b>	<b>Value Iteration Algorithm</b>	<b>10</b>
<b>7</b>	<b>Policy Iteration Algorithm</b>	<b>11</b>
<b>8</b>	<b>Comparison of two algorithms</b>	<b>11</b>
<b>9</b>	<b>Correctness of Algorithm</b>	<b>13</b>

---

# Math Foundation

## SECTION 1

### Basic Concept

---

#### SUBSECTION 1.1

#### Agent-Environment interaction cycle

---

The basic process in RL is called Agent-Environment interaction cycle. The interaction can go on for several cycles. Each cycle is called a time-step. Agent-Environment interaction cycle includes steps which are listed below:

1. Environment receives agent actions
2. Depending on current state and the agent chosen action, environment transitions to a new state
3. The new state and reward will pass to the agent, at most time these signals are passed through a filter because agent is forbidden to perceive the true environment.
4. According to the new state and reward, agent do the next action. <sup>1</sup>
5. **Deep Reinforcement Learning**

<sup>1</sup> Agent can do actions like learning or randomly choose next action

#### SUBSECTION 1.2

#### Episodic task, Continuing task and Return

---

Tasks that agent is trying to solve have a natural ending, are *episodic tasks*. Tasks that don't, are *continuing tasks*. A *time-step* limit is often added to *continuing tasks*.

An **episode** is a sequence of time-steps from the beginning to the end of an **episodic task**.<sup>2</sup> Agent may take several episodes to solve the *episodic task*.<sup>3</sup>

**Return** is the sum of all rewards that agent gets in an episode.

<sup>2</sup> like playing a game and in the end you can win or lose

<sup>3</sup> you should play again and again to win the game.

#### SUBSECTION 1.3

#### States

---

A **state** is a specific configuration of the problem (or environment). The set of all possible states is defined as the **state space**.

**State space** is a set of sets. We can treat it as two parts: the inner set and the outer set. The inner set must be finite, the elements of inner set are variables compose the states (variables do not have to be all inclusive). Elements in outer set represent possible states. Since the task can be a *continuing tasks*, the outer set may be infinite.

#### Independence

States must satisfy the **independence**. It means that the representation of states does not have to include all variables which compose the states, however, it must **contain all the variables necessary to make it independent of all other states**.

## MDPs vs POMDPs

**MDPs** is Markov Decision Process and **POMDPs** is Partially observable Markov Decision Process.

1. MDPs: the observation is the same as the states.
2. POMDPs (*more general*): the observation is part of the states

## Markov Property

Current RL or DRL systems are designed to take advantage of the Markov property. Variables (or the inner sets) you feed to the agent must **hold Markov assumption as tightly as possible**.

**Definition 1** (Markov assumption) The probability of the next state, given the current state and action, is dependent of the history of interactions.

$$p(S_{t+0}|S_t, A_t) = p(S_{t+1}|S_t, A_t, S_{t-1}, A_{t-1}, \dots) \quad (1.1)$$

## States in MDPs

The set of all states in the MDP is denoted as  $S^+$ . Two unique states in MDPs are starting state and **Terminal State**.<sup>4</sup>

A RL system can have only one terminal state, and also can have multiple terminal states.

SUBSECTION 1.4

## Actions

**Actions** are mechanisms to influence states. MDPs make available a set of actions  $A(s)$  that depends on states.  $A(s)$  is a function taking a state as its argument, and return some available actions base on the state.

The Action space with States space can be finite or infinite, but one action in action space must have finite variables to describe itself.

Unlike the number of variables compose the state, number of variables compose actions may not be constant because actions are depend on the state.

The environment makes all available actions known by agents in advance. Agent can select actions deterministically (from a look up table) or stochastically (from a per-state probability distribution).

SUBSECTION 1.5

## Transition Function

When agent takes one action, environment must respond to the action. The respond is transitioning current state to next state.

The way of changing is referred to as the **state-transition probabilities**, or for simplicity, **transition function**. According to its definition, transition function should have three parameters, current state  $s$ , action  $a$  taking at  $s$  and the next state  $s'$ . Transition function maps the tuple:  $(s, a, s')$  to a probability.

Transition function describes a probability distribution determining how RL system will evolve from current state to next state (with what probability change to a next specific state). As a probability distribution, if we sum or integrate over all next state

<sup>4</sup> Terminal state is one unique state in the set of states in a MDP. It must have all available actions transitioning, with probability 0, to itself, and these transitions must not provide any reward (reward equals 0).

$s'$ , the result must be one.

$$P(s'|s, a) = P(S_t = s' | S_{t-2} = s, A_{t-1} = a)$$

$$\sum_{s \in S} P(s'|s, a) = 0, \forall s \in S, a \in A(s)$$

#### SUBSECTION 1.6

### Rewards

The **reward function** maps a transition tuple:  $s, s', a$  to a scalar.

While the reward function can be represented by  $R(s, s', a)$ , which is explicit, it also can be represented by  $R(s, s', a)$  or even  $R(s)$ , depending on needs.

We can compute the marginalization over next state  $s'$ , to obtain  $R(s, a)$ , and the marginalization over action  $a$ , to obtain  $R(s)$ . For a stochastic process, reward function is defined by expectation.

$$r(s, a) = E[R_t | S_{t-2} = s, A_{t-1} = a]$$

$$r(s, a, s') = E[R_t | S_{t-2} = s, A_{t-1} = a, S_t = s']$$

#### SUBSECTION 1.7

### Horizon

A time-step, also referred to as **epoch, cycle, iteration or even interaction**, is a global clock syncing all parties and discrete time. Having a clock gives rise to a couple of types of task, *episodic task* and *continuing task*.

**Planning horizon** is agent's perspective that defines the *episodic task* and *continuing task*.

**Finite Horizon** is a planning horizon in which agent knows task will terminate in a finite number of steps. If the finite number of steps equals one, we call this planning horizon *greedy horizon*.

**Infinite Horizon** is the other planning horizon in which agent does not have pre-determined time step limit.

Although agent plans for infinite, interaction may be stopped by environment (common practice is adding an artificial terminal state). We refer the sequence of consecutive time steps from the beginning to the end of an episodic task as an **episode, trial, period, or stage**.<sup>5</sup>

In infinite horizon, episode also has its meaning that is a collection contains all interactions between an initial and a terminal state<sup>6</sup>.

<sup>5</sup> the same as the definition in section 1.2

<sup>6</sup> We can add an artificial terminal state

#### SUBSECTION 1.8

### Discount Factor

The **discount factor** adjust the importance of rewards, it helps reduce the degree to which future rewards affect our value function estimates and stabilize learning. Because that the future is uncertain, and the further we look into the future, the more stochastic we accumulate and the more variance our value function estimates will have.

Discount factor is part of MDP, not the agent. It is also a hyper parameter.

The **return**  $G_t$  is sum of all rewards obtained during an episode with discount factor and can be defined as below:

$$G_t = R_{t+0} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \quad (1.2)$$

The **recursive definition of return** is:

$$G_t = R_{t+0} + \gamma G_{t+1} \quad (1.3)$$

# Balancing immediate and long-term goals

This is about the basic algorithm for solving MDPs. It is a technique called **dynamic programming**, includes **Value Iteration (VI)** and **Policy Iteration (PI)**.

Although VI and PI require full access to the MDP, like knowing the dynamics of the environment, which is in some cheating way, they are the foundations from which virtually every other RL algorithm originates.

## SECTION 2

### Policy

Solid plan can not account for stochastic environment. What agent needs to come up with is called **policy**. Comparing with a plan, **policy** has a action for all possible states (only states in the plan have a action).

**Policies** can be stochastic (return action-probability distribution) or deterministic (return single actions for a given state). A policy is a function that prescribes actions to take for a given non-terminal state. Now, one immediate question is: *how good is a policy?*

## SECTION 3

### State-value function: What to expect from here?

If we are given a policy and the MDP, we should be able to calculate the expected return staring from every single state. That is, a number should be put into every state in the MDP for a given policy. The number must tell agent whether a state is better than another and precisely how much better it is.

Remember that the value of states depends on the policy which the agent follows. We define the value of states as the expectation of returns that the agent would get if it follows the policy  $\pi$ .

**Definition 2** (State-value function)

$$V_{\pi}(s) = E[G_t | S_t = s] \quad (3.1)$$

Considering the recursive definition of returns, the value of states can be rewrite as:

$$V_{\pi}(s) = E[R_{t+1} + \gamma G_{t+1} | S_t = s] \quad (3.2)$$

and this will derive famous Bellman Equation.

$$V_{\pi}(s) = \sum_a \pi(a|s) [\sum_r r P(r|s, a) + \gamma \sum_{s'} V_{\pi}(s') P(s'|s, a)] \quad \forall s, s' \in S \quad (3.3)$$

We can use margin probability to rewrite equation above, because

$$P(r|s, a) = \sum_{s'} P(r, s'|s, a) \quad P(s'|s, a) = \sum_r P(r, s'|s, a) \quad (3.4)$$

Then

$$V_\pi(s) = \sum_a \pi(a|s) \sum_r \sum_{s'} P(r, s'|s, a) [r + \gamma V_\pi(s')] \quad \forall s, s' \in S \quad (3.5)$$

**Definition 3** (Bellman Equation) We start from here:

$$V_\pi(s) = E[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

According to expectation definition and properties:

$$V_\pi(s) = E[R_{t+1} | S_t = s] + \gamma E[G_{t+1} | S_t = s]$$

The first item on equation right side can be expand as (according to Law of Total Probability<sup>7</sup>):

$$\begin{aligned} E[R_{t+1} | S_t = s] &= \sum_r r P(r|s) \\ &= \sum_r r \left( \sum_a \pi(a|s) P(r|s, a) \right) \end{aligned} \quad (3.6)$$

For simplicity, we write  $S_t = s, A_t = a$  as  $s, a$ . The second item on equation right side can be expand as (according to Properties of Expectation<sup>8</sup>)

$$\begin{aligned} E[G_{t+1} | S_t = s] &= \sum_{s'} E[G_{t+1} | S_t = s, S_{t+1} = s'] P(S_{t+1} = s' | S_t = s) \\ &= \sum_{s'} E[G_{t+1} | S_{t+1} = s'] P(s'|s) \quad (\text{Markov Property}) \\ &= \sum_{s'} V_\pi(s') \sum_a \pi(a|s) P(s'|s, a) \quad (\text{Law of Total Probability}) \end{aligned} \quad (3.7)$$

We combine equation (3.6) and (3.7), then get Bellman Equation.

$$V_\pi(s) = \sum_a \pi(a|s) \left[ \sum_r r P(r|s, a) + \gamma \sum_{s'} V_\pi(s') P(s'|s, a) \right] \quad \forall s, s' \in S$$

<sup>7</sup> If  $B_1, B_2, \dots$  form a partition of the sample space  $S$ , then we can calculate the probability of event  $A$  as:

$$P(A) = \sum_{B_i} P(A|B_i) P(B_i)$$

<sup>8</sup>

$$E[X] = \sum_y E[X|Y = y] P(y)$$

$$E[X] = E[E[X|Y]]$$

**PROOF** (Properties of Expectation)

$$\begin{aligned} E[X] &= \sum_x x P(x) \\ &= \sum_x \sum_y x P(x|Y = y) P(y) \\ &= \sum_y P(y) \sum_x x P(x|Y = y) \\ &= \sum_y E[X|Y = y] P(y) \end{aligned}$$

For  $E[X] = E[E[X|Y]]$ , we can use above equation:

$$\begin{aligned} E[X] &= \sum_y E[X|Y = y] P(y) \\ &= E[E[X|Y]] \quad (\text{Definition of Expectation}) \end{aligned}$$



|

□

## SECTION 4

## Action-value function: What should I expect from here if I do this action?

---

Another important question is the value of taking action  $a$  in a state  $s$ . By comparing **two actions under the same policy**, we can use **action-value function** to select a better one and then improve our policy.

**Action-value function** can be defined as:

Definition 4

(Action-value function)

$$q_{\pi}(s, a) = E[G_t | S_t = s, A_t = a] \quad (4.1)$$

We expand right side (use recursive definition of Return and Properties of Expectation):

$$\begin{aligned} E[G_t | S_t = s, A_t = a] &= \sum_r rP(r|s, a) + \gamma \sum_{s'} E[G_{t+1} | s, a, s']P(s'|s, a) \\ &= \sum_r rP(r|s, a) + \gamma \sum_{s'} E[G_{t+1} | s']P(s'|s, a) \\ &= \sum_r rP(r|s, a) + \gamma \sum_{s'} V_{\pi}(s')P(s'|s, a) \end{aligned}$$

By comparing with eq(3.5), we have

$$q_{\pi}(s, a) = \sum_r rP(r|s, a) + \gamma \sum_{s'} V_{\pi}(s')P(s'|s, a) \quad (4.2)$$

and

$$V_{\pi}(s) = \sum_a \pi(a|s)q_{\pi}(s, a) \quad (4.3)$$

## SECTION 5

## Optimal Policy

---

An optimal policy is an policy that for every state  $s$  can obtain expectation returns **greater than or equal to any other policy**. Now, an **Optimal State-value function** of state  $s$  and **Optimal action-value function** in state  $s$  of action  $a$  can be defined as:

$$\begin{aligned} V_*(s) &= \max_{\pi} V_{\pi}(s) \\ q_*(s, a) &= \max_{\pi} q_{\pi}(s, a) \end{aligned}$$

According to eq(3.5), the Optimal State-value function can be defined as:

**Definition 5** (Optimal State-value function)

$$\begin{aligned} V_*(s) &= \max_{\pi} V_{\pi}(s) \\ &= \max_{\pi} \sum_a \pi(a|s) \sum_{s'} \sum_r P(s', r|s, a) [r + \gamma V_{\pi}(s')] \end{aligned} \quad (5.1)$$

According to equation above, Optimal State-value function of state is determined by selecting an optimal policy that makes  $\sum_a \pi(a|s) \sum_{s'} \sum_r P(s', r|s, a) [r + \gamma V_{\pi}(s')]$  maximum (this is the meaning of  $\max_{\pi}$ ). Because  $\pi$  is an optimal policy,  $V_{\pi}(s') = V_*(s')$ . Then we should **clarify what action should be in an optimal policy**. In any states  $s$ , there must exist an action that maximizes  $\sum_{s'} \sum_r P(s', r|s, a) [r + \gamma V_*(s')]$ . Since  $\sum_a \pi(a|s) = 1$ , we have:

$$\begin{aligned} & \sum_a \pi(a|s) \sum_{s'} \sum_r P(s', r|s, a) [r + \gamma V_*(s')] \\ & \leq \sum_a \pi(a|s) \max_a \sum_{s'} \sum_r P(s', r|s, a) [r + \gamma V_*(s')] \\ & = \max_a \sum_{s'} \sum_r P(s', r|s, a) [r + \gamma V_*(s')] \end{aligned}$$

Now, we have the optimal policy  $\pi^*(a|s)$ ,

$$\pi^*(a|s) = \begin{cases} 1 & \text{if } a = a^*, \\ 0 & \text{if } a \neq a^*. \end{cases} \quad (5.2)$$

here  $a^* = \arg \max_a \sum_{s'} \sum_r P(s', r|s, a) [r + \gamma V_*(s')]$ . So we rewrite eq (5.1) as below:

$$V_*(s) = \max_a \sum_{s'} \sum_r P(s', r|s, a) [r + \gamma V_*(s')] \quad (5.3)$$

And according to eq.(4.2) and (4.3)

$$\begin{aligned} q_{\pi}(s, a) &= \sum_r r P(r|s, a) + \gamma \sum_{s'} V_{\pi}(s') P(s'|s, a) \\ &= \sum_r r P(r|s, a) + \gamma \sum_{s'} \sum_{a'} \pi(a'|s') q_{\pi}(s', a') P(s'|s, a) \\ &= \sum_r \sum_{s'} r P(r, s'|s, a) + \gamma \sum_r \sum_{s'} \left[ \sum_{a'} \pi(a'|s') q_{\pi}(s', a') P(r, s'|s, a) \right] \\ &= \sum_r \sum_{s'} P(r, s'|s, a) \left[ r + \gamma \sum_{a'} \pi(a'|s') q_{\pi}(s', a') \right] \end{aligned} \quad (5.4)$$

So Optimal Action-value function can be defined as:

**Definition 6** (Optimal Action-value function)

$$\begin{aligned}
 q_*(s, a) &= \max_{\pi} q_{\pi}(s, a) \\
 &= \max_{\pi} \sum_r \sum_{s'} P(r, s' | s, a) \left[ r + \gamma \sum_{a'} \pi(a' | s') q_{\pi}(s', a') \right] \\
 &= \sum_r \sum_{s'} P(r, s' | s, a) \left[ r + \gamma \max_{a'} q_*(s', a') \right] \tag{5.5}
 \end{aligned}$$

So far, we have defined the equation for Optimal State-value function and Optimal Action-value function. They are also the objectives we are after. We can start exploring the methods for finding these objectives.

# Two algorithm to find objects

The key point of both algorithms is the iterative method for solving Bellman optimality equation. In particular, the algorithm is:

$$\begin{aligned}
 & \text{iteratively compute} \\
 & V_{k+1}(s) = \max_{\pi} \sum_a \pi(a|s) \sum_{s'} \sum_r P(s', r|s, a) [r + \gamma V_k(s')] \\
 & = \max_{\pi} \sum_a \pi(a|s) q_k(s, a)
 \end{aligned} \tag{5.6}$$

In the end of this part, we will demonstrate why the iterative method can solve optimal Bellman Equation and find the object.

## SECTION 6

### Value Iteration Algorithm

Accroding to the eq.(5.6), we immeidately get the Value Iteration Algorithm which is listed below:

---

#### Algorithm 1: Value Iteration Algorithm

---

**Input:** The probability models  $P(r|s, a)$  and  $P(s'|s, a)$  for all  $(s, a)$ . Initial guess  $v_0$ .

**Output:** The optimal state value and optimal policy.

```

1 while  $V_k$  has not converged in the sense that  $\|V_k - V_{k-1}\|$  is greater
   than a predefined threshold do
2   foreach state  $s \in S$  do
3     foreach  $a \in A(s)$  do
4        $q_k(s, a) \leftarrow \sum_r P(r|s, a)r + \gamma \sum_{s'} P(s'|s, a)V_k(s')$ 
5     end
6      $a_k^*(s) = \arg \max_a q_k(s, a)$  // Maximum action value
7      $\pi_{k+1}(a|s) = 1$  if  $a = a_k^*(s)$ , else  $\pi_{k+1}(a|s) = 0$  // Policy updates
8      $V_{k+1}(s) \leftarrow \max_a q_k(s, a)$  // Value updates
9   end
10   $k \leftarrow k + 1$ 
11 end

```

---

In summary, the above steps can be illustrated as :

$V_k(s) \rightarrow q_k(s, a) \rightarrow$  greedly update policy  $\pi_{k+1}(a|s) \rightarrow$  new value  $V_{k+1}(s) = \max_a q_k(s, a)$

What should be noticed is that in Value Iteration algorithm  $V_{k+1}(s)$  is not the value function because it can not guarantee<sup>9</sup>

<sup>9</sup> In this equation, the index of  $q(s, a)$  is  $k$  not  $k + 1$ .

$$V_{k+1}(s) = \max_{\pi} \sum_a \pi(a|s) q_{k+1}(s, a)$$

If we want  $V_{k+1}(s)$  to be a value function, in the value updates step of Value Iteration algorithm, we should use an iterative procedure in place of the one step computation to get the value function under the updated policy  $\pi_{k+1}$ .

**Why?** Although  $V_{k+1}(s)$  is merely an intermediate value generated by the algorithm, the iterative algorithm can guarantee that  $V_{k+1}$  converge to the optimal state value function.

## SECTION 7

## Policy Iteration Algorithm

Not like Value Iteration algorithm that directly solve the Bellman optimality equation, **Policy Iteration** algorithm has an intimate relationship with **Value Iteration algorithm**.

---

**Algorithm 2:** Policy Iteration Alogrithm

---

```

Input: T
1 the system model,  $P(r|s, a)$  and  $P(s'|s, a)$  for all  $(s, a)$ , is known. Initial guess
    $\pi_0$ . Output: S
2 earch for the optimal state value and an optimal policy.
3 while  $V_{\pi_k}$  has not converged, for the  $k$ th iteration do
4   Initialization: randomly guess  $V_{\pi_k}^0$ 
   /* Iteratively find the approximation of the true state value
   function */
5   while  $V_{\pi_k}^j$  has not converged, for the  $j$ th iteration do
6     foreach state  $s \in S$  do
7        $V_{\pi_k}^{j+1} \leftarrow \sum_a \pi(a|s) \sum_r \sum_{s'} [P(r|s, a) + \gamma P(s'|s, a) V_{\pi_k}^j]$ 
8     end
9   end
   /* Policy improvement */
10  foreach state  $s \in S$  do
11    foreach action  $a \in A(s)$  do
12       $q_{\pi_k}(s, a) \leftarrow \sum_r P(r|s, a)r + \gamma \sum_{s'} P(s'|s, a) V_{\pi_k}(s')$ 
13    end
14     $a_k^*(s) = \arg \max_a q_{\pi_k}(s, a)$ 
15     $\pi_{k+1}(a|s) = 1$  if  $a = a_k^*(s)$ , else  $\pi_{k+1}(a|s) = 0$  // Policy updates
16  end
17   $k \leftarrow k + 1$ 
18 end

```

---

The embedded iterative algorithm of policy evaluation requires an infinite steps that is  $j \rightarrow \infty$  to converge to the true state value function. But it is impossible. We actually get an approximation of the true state value function but this does not cause any problems.<sup>10</sup>

## SECTION 8

## Comparison of two algorithms

Table below clearly illustrates the steps of two algorithms. We can see they are similar (if the initial guess of  $V_0$  in Value Iteration algorithm is replaced with  $V_{\pi_0}$  which is iteratively computed in Policy Iteration algorithm ). What different between Value

<sup>10</sup>  $V_{k+1}(s)$  in Value Iteration algorithm is more far from the true state value than that in policy evaluation, but it also can find the solution for optimal Bellman Equation.

**Table 1.** Comparison of two algorithm

	Policy iteration algorithm	Value iteration algorithm	Comments
1) Policy	$\pi_0$	N/A	
2) Value	$V_{\pi_0}$ ( <b>iteratively compute under <math>\pi_0</math></b> )	$V_0 = V_{\pi_0}$	replace random initial guess $V_0$ in Value iteration with $V_{\pi_0}$
3) Policy	$\pi_1 = \arg \max_{\pi} (r_{\pi} + \gamma P_{\pi} V_{\pi_0})$	$\pi_1 = \arg \max_{\pi} (r_{\pi} + \gamma P_{\pi} V_{\pi_0})$	two operations are the same
4) Value	$V_{\pi_1} = r_{\pi_1} + \gamma P_{\pi_1} V_{\pi_1}$	$V_1 = r_{\pi_1} + \gamma P_{\pi_1} V_{\pi_0}$	$V_{\pi_1}$ is iteratively computed whereas $V_1$ is just an one-step assignment
$\vdots$	$\vdots$	$\vdots$	

Iteration and Policy Iteration is in value compute, the  $V_{\pi_j}$  in Policy iteration is an approximation of Value State function but the  $V_j$  is just an intermediate result.

If we explicitly write the whole iterative procedure, we have

$$V_{\pi_1}^{(0)} = V_0$$

$$\text{Value iteration} \leftarrow V_1 \leftarrow V_{\pi_1}^{(1)} = r_{\pi_1} + \gamma P_{\pi_1} V_{\pi_1}^{(0)}$$

$$V_{\pi_1}^{(2)} = r_{\pi_1} + \gamma P_{\pi_1} V_{\pi_1}^{(1)}$$

$$\vdots$$

$$\text{truncated policy iteration} \leftarrow \bar{V}_1 \leftarrow V_{\pi_1}^{(j)} = r_{\pi_1} + \gamma P_{\pi_1} V_{\pi_1}^{(j-1)}$$

$$\vdots$$

$$\text{Policy iteration} \leftarrow V_{\pi_1} \leftarrow V_{\pi_1}^{(\infty)} = r_{\pi_1} + \gamma P_{\pi_1} V_{\pi_1}^{(\infty)}$$

and can have an observation from above process: if we truncated policy iteration, for example, limit the maximum number of iteration, we can get a new algorithm which is a middle ground between two extremes. It is called **Truncated Policy Iteration Algorithm**.

**Algorithm 3:** Truncated Policy Iteration algorithm

---

**Input:** T  
 1 he probability models  $P(r|s, a)$  and  $P(s'|s, a)$  for all  $(s, a)$  are known. Initial guess  $\pi_0$ . **Output:** T  
 2 he optimal State Value function and Policy.  
 3 **while**  $V_k$  (**not**  $V_{\pi_k}$ ) **has not converged, for the kth iteration do**  
 4     Select the initial value as  $V_k^{(0)} = V_{k-1}$   
       /\* truncate infinite j to a fixed number  $j_{truncate}$  \*/  
 5     **while**  $j < j_{truncate}$  **do**  
 6          $V_k^{(j)} = \sum_a \pi(a|s) \sum_r \sum_{s'} [P(r|s, a)r + \gamma P(s'|s, a)V_k^{(j-1)}(s')]$   
 7     **end**  
 8     Set  $V_k = V_k^{j_{truncate}}$   
       /\* Policy improvement \*/  
 9     **foreach** **state**  $s \in S$  **do**  
 10         **foreach** **action**  $a \in A(s)$  **do**  
 11              $q_k(s, a) = \sum_r P(r|s, a)r + \gamma \sum_{s'} P(s'|s, a)V_k(s')$   
 12         **end**  
 13          $a_k^*(s) = \arg \max_a q_k(s, a)$   
 14          $\pi_{k+1}(a|s) = 1$  if  $a = a_k^*$ , else  $\pi_{k+1}(a|s) = 0$  // Policy updates  
 15     **end**  
 16 **end**

---

## SECTION 9

**Correctness of Algorithm**

Nomatter Value Iteration or Policy Iteration, iterative method(5.6) for solving optimal Bellman Equation is the key point. The correctness of it is assured by **Contract Mapping Theorem**.

**Definition 7** (Contract Mapping) The function  $f(x)$  is a **contract mapping** if there exists  $\gamma \in (0, 1)$  such that

$$\|f(x_1) - f(x_2)\| \leq \gamma \|x_1 - x_2\| \quad (9.1)$$

for any  $x_1, x_2 \in R^d$ .

**Theorem 1** (Contract mapping theorem) For any equation that has the form  $x = f(x)$ , where  $x$  and  $f(x)$  are real vectors, if  $f$  is a contract mapping, then the following properties hold:

1. Existence: there exists a fixed point  $x^*$  satisfying  $f(x^*) = x^*$ .
2. Uniqueness: the fixed point is unique.
3. Algorithm: the following iterative algorithm

$$x_{k+1} = f(x_k) \quad (9.2)$$

where  $(k = 0, 1, \dots)$  will converge to  $(x^*)$  when  $(k \rightarrow \infty)$

Before proving this theorem, we give the definition of **fixed point** and **Cauchy sequence**.

**Definition 8** (Fixed point) Consider a function  $f(x)$ , where  $x \in R^d$  and  $f : R^d \rightarrow R^d$ . A point is  $x^*$  is a fixed point if  $f(x^*) = x^*$

**Definition 9** (Cauchy sequence) A sequence  $x_1, x_2, \dots \in R$  is called **Cauchy Sequence** if for  $\forall \epsilon > 0$ , there exists an integer  $N$  such that  $\|x_m - x_n\| < \epsilon$ , for all  $m, n > N$ .  
**Cauchy Sequence is guaranteed to converge to a limit.**

### PROOF 1. Convergent

For any  $m > n$ ,

$$\begin{aligned}\|x_m - x_n\| &= \|x_m - x_{m-1} + x_{m-1} - x_{m-2} + x_{m-2} - \dots + x_{n+1} - x_n\| \\ &\leq \|x_m - x_{m-1}\| + \|x_{m-1} - x_{m-2}\| + \dots + \|x_{n+1} - x_n\|\end{aligned}$$

and  $f$  is a contract mapping, we have

$$\begin{aligned}\|x_{k+1} - x_k\| &= \|f(x_k) - f(x_{k-1})\| \\ &\leq \gamma \|x_k - x_{k-1}\| \\ &= \gamma \|f(x_{k-1}) - f(x_{k-2})\| \\ &\leq \gamma^2 \|x_{k-1} - x_{k-2}\| \\ &\dots \\ &\leq \gamma^k \|x_1 - x_0\|\end{aligned}$$

So

$$\begin{aligned}\|x_m - x_n\| &\leq \gamma^{m-1} \|x_1 - x_0\| + \gamma^{m-2} \|x_1 - x_0\| + \dots + \gamma^n \|x_1 - x_0\| \\ &= \frac{\gamma^n (1 - \gamma^{m-n})}{1 - \gamma} \\ &\leq \frac{\gamma^n}{1 - \gamma} \|x_1 - x_0\|\end{aligned}\tag{9.3}$$

Now, for any  $\epsilon > 0$ , we can always find an integer  $N$  such that  $\|x_m - x_n\| < \epsilon$  for all  $m, n > N$ . So  $\{x_{k+1} = f(x_k)\}_{k=0}^\infty$  is a **Cauchy sequence** and hence converge to a limit point which is denoted as  $x^* = \lim_{k \rightarrow \infty} x_k$

### 2. The limit point is a fixed point

Because

$$\|f(x_k) - x_k\| = \|x_{k+1} - x_k\| \leq \gamma^k \|x_1 - x_0\|$$

$\|f(x_k) - x_k\|$  converges to zero exponentially. We get  $f(x^*) = x^*$

### 3. Uniqueness

Suppose there is another fixed point  $x'$  which satisfies  $f(x') = x'$ , then

$$\|x' - x^*\| = \|f(x') - f(x^*)\| \leq \gamma \|x' - x^*\| \quad (\text{definition of contract mapping}) \tag{9.4}$$

Since  $0 < \gamma < 1$ , this inequality holds if and only if  $\|x' - x^*\| = 0$  □

Contract mapping theorem tells us, if a mapping is a contract mapping, a Cauchy



sequence  $x_{k=0}^\infty$  can be generated by an iterative method ( $x_{k+1} = f(x_k)$ ) and its limit is the fixed point of that mapping.

Here we demonstrate the right-hand side of equation  $V_*(s)$  is a contract mapping. Firstly, we define  $r_\pi(s)$  and  $p_\pi(s'|s)$  based on eq.(3.5).

$$r_\pi(s) = \sum_a \pi(a|s) \sum_r r P(r|s, a)$$

$$P_\pi(s'|s) = \sum_a \pi(a|s) P(s'|s, a)$$

Suppose that the states are indexed as  $s_i$  with  $i = 1, 2, \dots, n$ , for state  $s_i$  in eq.(3.5)

$$V_\pi(s_i) = r_\pi(s_i) + \gamma \sum_{s_j} p_\pi(s_j|s_i) V_\pi(s_j)$$

Let  $V_\pi = [V_\pi(s_1), V_\pi(s_1), \dots, V_\pi(s_n)]^T \in \mathbb{R}^n$ ,  $r_\pi = [r_\pi(s_1), r_\pi(s_1), \dots, r_\pi(s_n)]^T \in \mathbb{R}^n$ , and  $P_\pi \in \mathbb{R}^{n \times n}$  with  $[P_\pi]_{ij} = P_\pi(s_i|s_j)$  Then we get the matrix form of eq.(3.5)

$$V_\pi = r_\pi + \gamma P_\pi V_\pi$$

Here is an example with four states:

$$\underbrace{\begin{bmatrix} V_\pi(s_1) \\ V_\pi(s_2) \\ V_\pi(s_3) \\ V_\pi(s_4) \end{bmatrix}}_{V_\pi} = \underbrace{\begin{bmatrix} r_\pi(s_1) \\ r_\pi(s_2) \\ r_\pi(s_3) \\ r_\pi(s_4) \end{bmatrix}}_{r_\pi} + \gamma \underbrace{\begin{bmatrix} P_\pi(s_1|s_1) & P_\pi(s_2|s_1) & P_\pi(s_3|s_1) & P_\pi(s_4|s_1) \\ P_\pi(s_1|s_2) & P_\pi(s_2|s_2) & P_\pi(s_3|s_2) & P_\pi(s_4|s_2) \\ P_\pi(s_1|s_3) & P_\pi(s_2|s_3) & P_\pi(s_3|s_3) & P_\pi(s_4|s_3) \\ P_\pi(s_1|s_4) & P_\pi(s_2|s_4) & P_\pi(s_3|s_4) & P_\pi(s_4|s_4) \end{bmatrix}}_{P_\pi} \underbrace{\begin{bmatrix} V_\pi(s_1) \\ V_\pi(s_2) \\ V_\pi(s_3) \\ V_\pi(s_4) \end{bmatrix}}_{V_\pi}$$

Since under different state we should have different optimal policy, the value of policy  $\pi$  is determined by the states. The matrix form of Bellman Optimality Equation is:

$$V^* = \max_{\pi} (r_\pi + \gamma P_\pi V^*) \quad (9.5)$$

Right now, we can see if taking right-hand side of eq.(9.5) as a function  $f$ , optimal value function ( $V^*$ ) will be the fixed point of  $f(V)$ . That is why we will demonstrate  $f$  is a contract mapping.

Considering given any two vectors  $V_1, V_2 \in \mathbb{R}^n$ , and  $\pi_1 = \arg \max_{\pi} (r_\pi + \gamma P_\pi V_1)$  and  $\pi_2 = \arg \max_{\pi} (r_\pi + \gamma P_\pi V_2)$  (notice that policy  $\pi$  is determined by states and that is why the mapping taking only  $V$  as its variable rather than  $V$  and  $\pi$ ). Then

$$f(V_1) = r_{\pi_1} + \gamma P_{\pi_1} V_1 \geq r_{\pi_2} + \gamma P_{\pi_2} V_1$$

$$f(V_2) = r_{\pi_2} + \gamma P_{\pi_2} V_2 \geq r_{\pi_1} + \gamma P_{\pi_1} V_2$$

where  $\geq$  is an elementwise comparison.

As a result,

$$\begin{aligned} f(V_1) - f(V_2) &\leq r_{\pi_1} + \gamma P_{\pi_1} V_1 - r_{\pi_1} + \gamma P_{\pi_1} V_2 \\ &= \gamma P_{\pi_1} (V_1 - V_2) \end{aligned}$$

$$\begin{aligned} f(V_2) - f(V_1) &\leq r_{\pi_2} + \gamma P_{\pi_2} V_2 - r_{\pi_2} + \gamma P_{\pi_2} V_1 \\ &= \gamma P_{\pi_2} (V_2 - V_1) \end{aligned}$$

Therefore,

$$\gamma P_{\pi_2} (V_1 - V_2) \leq f(V_1) - f(V_2) \leq \gamma P_{\pi_1} (V_1 - V_2)$$

Define  $z \doteq \max(|\gamma P_{\pi_1} (V_1 - V_2)|, |\gamma P_{\pi_2} (V_1 - V_2)|)$  and use infinite norm ( $\|\cdot\|_\infty \doteq \max_i |x_i|$ ), we get

$$\|f(V_1) - f(V_2)\|_\infty \leq \|z\|_\infty$$

and

$$z_i = \max(\gamma |P_{\pi_1}[i](V_1 - V_2)|, \gamma |P_{\pi_2}[i](V_1 - V_2)|)$$

where  $z_i$  is the  $i$ th element in  $z$  and  $P_{\pi_{1,2}}[i]$  is the  $i$ th row in matrix  $P_{\pi_{1,2}}$ . Because  $P_{\pi_{1,2}}[i]$  is a vector with all non-negative elements and the sum of the elements are equal to 1, it follows that

$$|P_{\pi_1}[i](V_1 - V_2)| = P_{\pi_1}[i](V_1 - V_2) \leq \|V_1 - V_2\|_\infty,$$

Similarly,

$$|P_{\pi_2}[i](V_1 - V_2)| = P_{\pi_2}[i](V_1 - V_2) \leq \|V_1 - V_2\|_\infty.$$

So

$$\|z\|_\infty = \max_i |z_i| \leq \gamma \|V_1 - V_2\|_\infty$$

Finally we conclude that  $f(V)$  is a contract mapping because

$$\|f(V_1) - f(V_2)\|_\infty \leq \gamma \|V_1 - V_2\|_\infty \leq \|V_1 - V_2\|_\infty$$

According to **Contract mapping theorem**, we can easily get the theorem:

**Theorem 2** For the Bellman Optimality Equation  $V = f(V) = \max_\pi (r_\pi + \gamma P_\pi V)$ , there always exists a unique solution  $V^*$ , which can be iteratively solved by

$$V_{k+1} = \max_\pi (r_\pi + \gamma P_\pi V_k) \quad (9.6)$$

The value of  $V_k$  is converged exponentially to  $V^*$  as  $k \rightarrow \infty$  given any initial state  $V_0$

For the two algorithm Value Iteration and Policy Iteration, Value Iteration directly uses iterative method to solve BOE (Bellman Optimality Equation) whereas Policy Iteration just improve policy. So we can confirm that Value Iteration can get the optimal state value function and optimal policy<sup>11</sup>, but for Policy Iteration, we need to prove that it eventually find an optimal policy and optimal state value function.

**Lemma 1** (Policy improvement) In the policy improvement step,  $\pi_{k+1}$  is better than  $\pi_k$ . That is, if  $\pi_{k+1} = \arg \max_\pi (r_\pi + \gamma P_\pi V_{\pi_k})$ , then  $V_{\pi_{k+1}} \geq V_{\pi_k}$ .

<sup>11</sup> Solving  $\pi^*$  can be easily obtained by  $\pi^* = \arg \max_\pi (r_\pi + \gamma P_\pi V^*)$  once optimal State value function is obtained.

PROOF Since  $V_{\pi_{k+1}}$  and  $V_{\pi_k}$  are state value function, they satisfy the Bellman Equation:

$$\begin{aligned} V_{\pi_{k+1}} &= r_{\pi_{k+1}} + \gamma P_{\pi_{k+1}} V_{\pi_{k+1}} \\ V_{\pi_k} &= r_{\pi_k} + \gamma P_{\pi_k} V_{\pi_k} \end{aligned}$$

Because  $\pi_{k+1} = \arg \max_{\pi} (r_{\pi} + \gamma P_{\pi} V_{\pi_k})$ , we know that,

$$r_{\pi_{k+1}} + \gamma P_{\pi_{k+1}} V_{\pi_k} \geq r_{\pi_k} + \gamma P_{\pi_k} V_{\pi_k}$$

It then follows that,

$$\begin{aligned} V_{\pi_k} - V_{\pi_{k+1}} &= (r_{\pi_k} + \gamma P_{\pi_k} V_{\pi_k}) - (r_{\pi_{k+1}} + \gamma P_{\pi_{k+1}} V_{\pi_{k+1}}) \\ &\leq (r_{\pi_{k+1}} + \gamma P_{\pi_{k+1}} V_{\pi_{k+1}}) - (r_{\pi_{k+1}} + \gamma P_{\pi_{k+1}} V_{\pi_k}) \\ &= \gamma P_{\pi_{k+1}} (V_{\pi_k} - V_{\pi_{k+1}}) \end{aligned}$$

Therefore,

$$\begin{aligned} V_{\pi_k} - V_{\pi_{k+1}} &\leq \gamma^2 P_{\pi_{k+1}}^2 (V_{\pi_k} - V_{\pi_{k+1}}) \leq \dots \\ &\leq \gamma^n P_{\pi_{k+1}}^n (V_{\pi_k} - V_{\pi_{k+1}}) \\ &\leq \lim_{n \rightarrow \infty} \gamma^n P_{\pi_{k+1}}^n (V_{\pi_k} - V_{\pi_{k+1}}) = 0 \end{aligned}$$

The limit is due to  $\gamma < 1$  and  $P_{\pi_{k+1}}^n$  is a non-negative stochastic matrix for any  $n$   $\square$

Before demonstrating that policy iteration algorithm can get optimal value function, we recall converge of sequence. For monotonic sequence, its convergence can be guaranteed by theorem below:

**Theorem 3** (Convergence of monotonic sequences) If the sequence  $\{x_k\}$  is non-increasing and bounded from below:

1. Non-increasing:  $x_{k+1} \leq x_k$
2. Similarly, if  $\{x_k\}$  is non-decreasing and have a upper bound, then the sequence is convergent

then  $x_k$  converges to a value, which is the limit of  $\{x_k\}$ , as  $k \rightarrow \infty$ .

Similarly, if  $\{x_k\}$  is non-decreasing and have a upper bound, then the sequence is convergent.

Considering the non-monotonic sequence, we also have a similiar convergence.

**Theorem 4** (Convergence of non-monotonic sequence) if the sequences is non-negative  $x_k \geq 0$  and satisfies:

$$\sum_{k=1}^{\infty} (x_{k+1} - x_k)^+ < \infty \quad (9.7)$$

then  $\{x_k\}$  converges as  $k \rightarrow \infty$ .

The definition of symbol  $(\cdot)^+$  is: For any  $z \in \mathbb{R}$ , define

$$z^+ \doteq \begin{cases} z & \text{if } z \geq 0, \\ 0 & \text{if } z < 0. \end{cases}$$

Similarly, the definition of symbol  $(\cdot)^-$  is:

$$z^- \doteq \begin{cases} z & \text{if } z \leq 0, \\ 0 & \text{if } z > 0. \end{cases}$$

PROOF To analyze the convergence of non-monotonic sequences  $\{x_k\}$ , we can use symbols above to rewrite it as:

$$\begin{aligned} x_k &= x_k - x_{k-1} + x_{k-1} - x_{k-2} + \cdots + x_2 - x_1 + x_1 \\ &= \sum_{i=1}^{k-1} (x_{i+1} - x_i) + x_1 \\ &\doteq S_k + x_1 \end{aligned}$$

Note  $S_k = \sum_{i=1}^{k-1} (x_{i+1} - x_i)$  can be decomposed as:

$$S_k = S_k^+ + S_k^-$$

where

$$\begin{aligned} S_k^+ &= \sum_{i=1}^{k-1} (x_{i+1} - x_i), \text{ for all } i \text{ which } x_{i+1} - x_i \geq 0 \\ S_k^- &= \sum_{i=1}^{k-1} (x_{i+1} - x_i), \text{ for all } i \text{ which } x_{i+1} - x_i < 0 \end{aligned}$$

Obviously, sequence  $\{S_k^+\}$  is non-decreasing whereas  $\{S_k^-\}$  is non-increasing.

Because  $x_k \geq 0$ , then  $S_k^- \geq -S_k^+ - x_1$ . So if  $\{S_k^+\}$  has a upper bound,  $\{S_k^-\}$  also has a lower bound, and vice versa.

Now the convergence of non-negative and non-monotonic sequences which satisfy the condition 9.7 can be easily proved. First, the condition 9.7 indicates that  $\{S_k^+\}$  has a upper bound. Since  $\{S_k^+\}$  is non-decreasing, the convergence of  $\{S_k^+\}$  immediately follows from theorem 1. Suppose that  $\{S_k^+\}$  converges to  $S_*^+$ .

Second, because  $\{S_k^+\}$  has a upper bound, the non-increasing sequence  $\{S_k^-\}$  also has a lower bound. We also immediately get the convergence of  $\{S_k^-\}$  based on theorem 1. Suppose that  $\{S_k^-\}$  converges to  $S_*^-$ .

Finally, because  $x_k = S_k^+ + S_k^- + x_1$ , the convergence of  $\{S_k^+\}$  and  $\{S_k^-\}$  implies that  $\{x_k\}$  converges to  $x_1 + S_*^+ + S_*^-$ .  $\square$

In Policy iteration, we obtain two sequence, one is the sequence of policies:

$$\{\pi_1, \pi_2, \dots, \pi_k, \dots\}$$

and the other is the sequence of state value function:

$$\{V_{\pi_1}, V_{\pi_2}, \dots, V_{\pi_k}, \dots\}.$$

Based on the Policy Improvement Lemma, we know that the sequence of state value function is a non-decreasing sequence and bounded by  $V_*$  and it follows from **Convergence of monotonic sequences** that  $V_{\pi_k}$  converges to a constant value  $V_\infty$ , when  $k \rightarrow \infty$ . The final step is to demonstrate  $V_\infty = V_*$ <sup>12</sup>.

<sup>12</sup> Convergent value may not equal to bound!

*Remark*

The idea of the proof is to show that the policy algorithm has a faster convergent speed than the value iteration. What the idea means is that as value iteration can converges to the optimal value  $V_*$ , with **squeeze theorem** ( $V_*$  is the lower bound of the sequence of state value function which is generated by policy iteration algorithm since policy iteration has a faster convergent speed than value iteration and according to the definition,  $V_*$  is the upper bound of the sequence), policy iteration is able to converge to  $V_*$ .

Here is the proof,

**PROOF**

In particular, to prove the convergence of  $\{V_{\pi_k}\}_{k=0}^{\infty}$ , we introduce another sequence  $\{V_k\}_{k=0}^{\infty}$  generated by

$$V_{k+1} = \max_{\pi_k} (r_{\pi_k} + \gamma P_{\pi_k} V_k) \quad (9.8)$$

In fact,  $\{V_k\}_{k=0}^{\infty}$  is exactly generated by value iteration. We use mathematical induction to demonstrate that the policy iteration can faster converge.

- 1) For  $k = 0$ , we can always find a value  $V_0$  that  $V_{\pi_0} \geq V_0$  under any initial guess  $\pi_0$
- 2) Suppose that for  $k \geq 0$ , we have  $V_{\pi_k} \geq V_k$
- 3) For  $k + 1$ ,

$$\begin{aligned} V_{\pi_{k+1}} - V_{k+1} &= (r_{\pi_{k+1}} + \gamma P_{\pi_{k+1}} V_{\pi_{k+1}}) - \max_{\pi_k} (r_{\pi_k} + \gamma P_{\pi_k} V_k) \\ &\geq (r_{\pi_{k+1}} + \gamma P_{\pi_{k+1}} V_{\pi_k}) - \max_{\pi_k} (r_{\pi_k} + \gamma P_{\pi_k} V_k) \quad (\text{Lemma 1 } V_{\pi_{k+1}} \geq V_{\pi_k}) \end{aligned}$$

Let  $\pi'_k = \max_{\pi_k} (r_{\pi_k} + \gamma P_{\pi_k} V_k)$ . We should notice that in policy improvement step of both value and policy iteration, they travel throughout every state and every action. That means the candidate policy set for the optimal policy in the  $k$ th iteration of policy iteration algorithm is as same as that in value iteration.

So if  $\pi_k^*$  which is equal to  $\pi_{k+1}$  (policy update step in Policy Iteration algorithm) is the optimal policy generated by the  $k$ th iteration of policy iteration algorithm and  $\pi'_k = \max_{\pi_k} (r_{\pi_k} + \gamma P_{\pi_k} V_k)$ , we have  $(r_{\pi_{k+1}} + \gamma P_{\pi_{k+1}} V_{\pi_k}) \geq (r_{\pi'_k} + \gamma P_{\pi'_k} V_{\pi_k})$  and,

$$\begin{aligned} V_{\pi_{k+1}} - V_{k+1} &\geq (r_{\pi'_k} + \gamma P_{\pi'_k} V_{\pi_k}) - (r_{\pi'_k} + \gamma P_{\pi'_k} V_k) \\ &= \gamma P_{\pi'_k} (V_{\pi_k} - V_k) \end{aligned}$$

Because  $V_{\pi_k} \geq V_k$  and  $P_{\pi'_k}$  is non-negative, we have  $V_{\pi_{k+1}} \geq V_{k+1}$ . Therefore we can show by induction that  $V_k \leq V_{\pi_k} \leq V^*$ . Since  $\{V_k\}_{k=0}^{\infty}$  converges to  $V^*$ ,  $V_{\pi_k}$  also converges to  $V^*$ .  $\square$

Now we have

**Theorem 5**

(Convergence of Policy Iteration) The state value sequence  $\{V_{\pi_k}\}_{k=0}^{\infty}$  generated by Policy Iteration converges to the optimal state value  $V_*$ . As a result, the policy sequence  $\{\pi_k\}_{k=0}^{\infty}$  converges to the optimal policy  $\pi^*$