

# Reinforecement Learning: Notes and Selected Exercises

Somnath Sikdar

November 12, 2020

# Contents

<b>1</b>	<b>Markov Decision Processes</b>	<b>2</b>
1.1	Introduction . . . . .	2
1.1.1	Policy and Value Functions . . . . .	3
1.1.2	Optimal Policies and Optimal Value Functions . . . . .	5

# Chapter 1

## Markov Decision Processes

### 1.1 Introduction

Reinforcement learning consists in teaching an agent how to behave using observations from its environment with the goal of maximizing future rewards that it receives in response to its behavior. The set-up consists of an *agent*, which is the learner or decision maker and the *environment*, which is the thing that the agent interacts with, comprising everything outside the agent. The agent and the environment interact in discrete time steps. At a given time step  $t$ , the agent takes an *action* and the environment responds by presenting a new state and a scalar feedback signal called the *reward*.

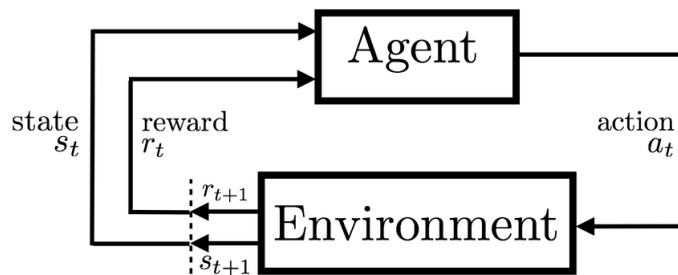


Figure 1.1: The agent and the environment

Problems in reinforcement learning are typically formulated in the formal setting of a Markov Decision Process (MDP). An MDP is a five-tuple  $\langle \mathcal{S}, \mathcal{A}, p, \mathcal{R}, \gamma \rangle$

- a set  $\mathcal{S}$  of states
- a set  $\mathcal{A}$  of actions
- a set  $\mathcal{R}$  of rewards
- a probability function  $p: \mathcal{S} \times \mathcal{R} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$
- a discount factor  $\gamma \in [0, 1]$

The probability function  $p$  defines the dynamics of the MDP. It tells us what the most likely next state and reward combination are given the current state and the action taken by the agent.

$$p(s', r \mid s, a) := \mathbf{Pr} \{S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a\}. \quad (1.1)$$

In the above equation, we assume that the probability does *not* depend on the time step  $t$ . Such an MDP is called *stationary*. Since this is a probability distribution, for all  $s \in \mathcal{S}$  and for all  $a \in \mathcal{A}$ ,

$$\sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r \mid s, a) = 1. \quad (1.2)$$

We also define a state-transition probability function  $p: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  defined as:

$$p(s' \mid s, a) := \sum_{r \in \mathcal{R}} p(s', r \mid s, a). \quad (1.3)$$

In response to the action taken by the agent, the environment provides it a scalar rewards at each time step. The total reward of the agent in time step  $t$  is defined as

$$G_t := \sum_{k=0}^{\infty} R_{t+k+1}. \quad (1.4)$$

With infinite time horizons, this reward could go to infinity. To prevent this, one uses a discounting factor  $\gamma \in [0, 1]$  and defines the total reward as

$$G_t := R_{t+1} + \gamma \cdot R_{t+2} + \gamma^2 \cdot R_{t+3} + \dots \quad (1.5)$$

### 1.1.1 Policy and Value Functions

A *policy*  $\pi$  is a mapping from states to probabilities of selecting each possible action.

$$\pi(a \mid s) := \mathbf{Pr} \{A_t = a \mid S_t = s\}. \quad (1.6)$$

Note that this probability distribution does not depend on the time step  $t$ .

**Exercise 1.1.** If the current state is  $S_t$ , and actions are selected according to a stochastic policy  $\pi$ , then what is the expectation of  $R_{t+1}$  in terms of  $\pi$  and the four-argument function  $p$ ?

**Solution.** Let us assume that the current state  $S_t = s$ . We may write  $\mathbf{E}_\pi [R_{t+1} \mid S_t = s]$  as:

$$\begin{aligned} \mathbf{E}_\pi [R_{t+1} \mid S_t = s] &= \sum_{r \in \mathcal{R}} r \cdot \mathbf{Pr}_\pi \{R_{t+1} = r \mid S_t = s\} \\ &= \sum_{r \in \mathcal{R}} r \cdot \sum_{a \in \mathcal{A}} \mathbf{Pr}_\pi \{R_{t+1} = r \mid S_t = s, A_t = a\} \cdot \mathbf{Pr}_\pi \{A_t = a \mid S_t = s\} \\ &= \sum_{r \in \mathcal{R}} r \cdot \sum_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} p(s', r \mid s, a) \cdot \pi(a \mid s) \end{aligned}$$

□

The *value function* of a state  $s$  under a policy  $\pi$ , denoted  $v_\pi(s)$ , is the expected return when starting in state  $s$  and following  $\pi$  thereafter. For an MDP, we may write the value function as:

$$v_\pi(s) := \mathbf{E}_\pi [G_t \mid S_t = s] = \mathbf{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k \cdot R_{t+k+1} \mid S_t = s \right]. \quad (1.7)$$

The function  $v_\pi(s)$  is called the *state-value function* of the policy  $\pi$ .

We also define the value of taking an action  $a$  in state  $s$  under a policy  $\pi$ , denoted  $q_\pi(s, a)$ , as the expected return when starting in state  $s$ , taking action  $a$  and following  $\pi$  thereafter.

$$q_\pi(s, a) := \mathbf{E}_\pi [G_t \mid S_t = s, A_t = a] = \mathbf{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k \cdot R_{t+k+1} \mid S_t = s, A_t = a \right]. \quad (1.8)$$

The function  $q_\pi(s, a)$  is the *action-value function* of the policy  $\pi$ .

**Exercise 1.2.** Give an equation for  $v_\pi$  in terms of  $q_\pi$  and  $\pi$ .

**Solution.** We may write  $v_\pi(s)$  as:

$$\begin{aligned} v_\pi(s) &:= \mathbf{E}_\pi [G_t \mid S_t = s] \\ &= \sum_{a \in \mathcal{A}} \mathbf{E}_\pi [G_t \mid S_t = s, A_t = a] \cdot \mathbf{Pr}_\pi \{A_t = a \mid S_t = s\} \\ &= \sum_{a \in \mathcal{A}} q_\pi(s, a) \pi(a \mid s). \end{aligned}$$

□

**Exercise 1.3.** Give an equation for  $q_\pi(s, a)$  in terms of  $v_\pi$  and the four parameter function  $p$ .

**Solution.**

$$\begin{aligned} q_\pi(s, a) &:= \mathbf{E}_\pi [R_{t+1} + \gamma \cdot R_{t+2} + \gamma^2 \cdot R_{t+3} + \cdots \mid S_t = s, A_t = a] \\ &= \mathbf{E}_\pi [R_{t+1} \mid S_t = s, A_t = a] + \gamma \cdot \mathbf{E}_\pi [G_{t+1} \mid S_t = s, A_t = a] \\ &= \sum_{r \in \mathcal{R}} r \cdot \mathbf{Pr}_\pi \{R_{t+1} = r \mid S_t = s, A_t = a\} + \\ &\quad \gamma \cdot \sum_{s' \in \mathcal{S}} \mathbf{E}_\pi [G_{t+1} \mid S_{t+1} = s', S_t = s, A_t = a] \cdot \mathbf{Pr}_\pi \{S_{t+1} = s' \mid S_t = s, A_t = a\} \\ &= \sum_{r \in \mathcal{R}} r \cdot \sum_{s' \in \mathcal{S}} p(s', r \mid s, a) + \gamma \cdot \sum_{s' \in \mathcal{S}} \mathbf{E}_\pi [G_{t+1} \mid S_{t+1} = s'] \cdot \sum_{r \in \mathcal{R}} p(s', r \mid s, a) \\ &= \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} [r + \gamma \cdot v_\pi(s')] p(s', r \mid s, a). \end{aligned}$$

Based on the last two exercises, we may write:

$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a \mid s) \cdot \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} [r + \gamma \cdot v_\pi(s')] p(s', r \mid s, a). \quad (1.9)$$

**Exercise 1.4.** Give an equation for  $q_\pi(s, a)$  in terms of the value  $q_\pi(s', a')$  of the next state-action pair  $(s', a')$ .

**Solution.** The probability that the next state-action pair is  $(s', a')$  given that the current state-action pair is  $(s, a)$  is:

$$\Pr_\pi \{S_{t+1} = s', A_{t+1} = a' | S_t = s, A_t = a\} = p(s'|s, a) \cdot \pi(a'|s').$$

Now,

$$\begin{aligned} q_\pi(s, a) &= \mathbf{E}_\pi [G_t | S_t = s, A_t = a] \\ &= \mathbf{E}_\pi [R_{t+1} | S_t = s, A_t = a] + \gamma \cdot \mathbf{E}_\pi [G_{t+1} | S_t = s, A_t = a] \\ &= \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} r \cdot p(s', r | s, a) + \\ &\quad \gamma \cdot \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} \mathbf{E}_\pi [G_{t+1} | S_{t+1} = s', A_{t+1} = a', S_t = s, A_t = a] \cdot p(s'|s, a) \cdot \pi(a'|s') \\ &= \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} r \cdot p(s', r | s, a) + \gamma \cdot \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} q_\pi(s', a') \cdot p(s'|s, a) \cdot \pi(a'|s') \end{aligned}$$

### 1.1.2 Optimal Policies and Optimal Value Functions

A policy  $\pi$  is said to be better than or equal to a policy  $\pi^*$  if for all states  $s \in \mathcal{S}$ ,  $v_\pi(s) \geq v_{\pi^*}(s)$ . An optimal policy  $\pi^*$  is one which satisfies the condition:

$$v_{\pi^*}(s) = \max_{\pi} v_\pi(s) \quad \forall s \in \mathcal{S} \quad (1.10)$$

In an MDP, one is guaranteed to have an optimal policy. Define  $v_\star := v_{\pi^*}$ . The following claim is easy to show.

**Lemma 1.1.** *Let  $\pi^*$  be an optimal policy. Then  $q_{\pi^*}(s, a) = \max_{\pi} q_\pi(s, a)$  for all  $s \in \mathcal{S}$  and all  $a \in \mathcal{A}$ .*

*Proof.* As shown in Exercise 1.3,  $\forall s \in \mathcal{S}$  and  $\forall a \in \mathcal{A}(s)$

$$\begin{aligned} q_{\pi^*}(s, a) &= \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} [r + \gamma \cdot v_\star(s')] p(s', r | s, a) \\ &= \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} \left[ r + \gamma \cdot \max_{\pi} v_\pi(s') \right] p(s', r | s, a) \\ &= \max_{\pi} \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} [r + \gamma \cdot v_\pi(s')] p(s', r | s, a) \\ &= \max_{\pi} q_\pi(s, a). \end{aligned}$$

□

Since  $v_*$  is the optimal state value function, it must equal the expected return for the best action from any given state. Hence,

$$\begin{aligned} v_*(s) &= \max_{a \in \mathcal{A}(s)} q_*(s, a) \\ &= \max_{a \in \mathcal{A}(s)} \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} [r + \gamma \cdot v_*(s')] p(s', r | s, a) \quad (\text{Exercise 1.3}). \end{aligned} \quad (1.11)$$

Similarly, we may start with the expression for  $q_*$  as in Exercise 1.3 and write:

$$\begin{aligned} q_*(s, a) &= \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} [r + \gamma \cdot v_*(s')] p(s', r | s, a) \\ &= \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} \left[ r + \gamma \cdot \max_{a' \in \mathcal{A}(s')} q_*(s', a') \right] p(s', r | s, a). \end{aligned} \quad (1.12)$$

Equations 1.11 and 1.12 are called the Bellman optimality equations for  $v_*$  and  $q_*$ , respectively.