

# Markov Decision Processes

Guillaume Barnier

Academic Year 2020-2021

## Contents

1	Markov Process	1
2	Markov Reward Process	2
3	Other useful tricks	4
4	Acknowledgments	4

## 1 Markov Process

**Definition 1.1** (Markov Property). Let  $S_t = (s_0, s_1, \dots)$  be a stochastic process evolving according to a transition dynamic  $P$ . This stochastic process satisfies the Markov property if

$$P(s_t | s_0, s_1, \dots, s_{t-1}) = P(s_t | s_{t-1}), \forall t \in \mathbb{N} \quad (1)$$

**Definition 1.2** (Markov Process). A Markov Process (MP) is a stochastic process that satisfies the Markov property

In a RL setting, we often make two additional assumptions:

- **Finite state space.** The state space of the Markov process is finite. This means that for the Markov process  $(s_0, s_1, \dots)$ , there is a state space  $S$  with  $|S| < \infty$ , such that for all realizations of the Markov process, we have  $s_t \in S$  for all  $t$ .
- **Stationary transition probability.** The transition probabilities are time independent:

$$P(s_p = s' | s_{p-1} = s) = P(s_q = s' | s_{q-1} = s), \forall (p, q) \quad (2)$$

A Markov process satisfying these assumptions is also sometimes called a Markov chain, although the precise definition of a Markov chain varies. With these assumptions, we can define characterize a Markov process with the following definition.

**Definition 1.3** (Markov Process/Markov Chain). A Markov Process is a tuple  $(S, P)$ , where

- $S$  is the finite state-space of the Markov process,  $|S| < \infty$
- $P$  is the state transition probability model where  $P_{ss'} = P[s_{t+1} = s' | s_t = s]$

**Lemma 1.1.**  $P(s_{t+n} | s_t = s) = P(s_n | s_0 = s)$  for all  $t$  and  $n$

*Proof.* I show this property by recursion on  $n$ :

- For  $n = 1$ ,  $P(s_{t+1} | s_t = s) = P(s_1 | s_0 = s)$  is true by the stationarity assumption
- I assume that  $P(s_{t+n} | s_t = s) = P(s_n | s_0 = s)$  is true for  $n$

- I show that  $P(s_{t+n+1}|s_t = s) = P(s_{n+1}|s_0 = s)$ :

$$P(s_{t+n+1}|s_t = s) = \sum_{s'} P(s_{t+n+1}, s_{t+n} = s' | s_t = s) \quad (3)$$

$$= \sum_{s'} P(s_{t+n+1}|s_t = s, s_{t+n} = s') P(s_{t+n} = s' | s_t = s) \quad (4)$$

$$= \sum_{s'} P(s_{t+n+1}|s_{t+n} = s') P(s_n = s' | s_0 = s) \quad (5)$$

$$= \sum_{s'} P(s_{n+1}|s_n = s') P(s_n = s' | s_0 = s) \quad (6)$$

$$= \sum_{s'} P(s_{n+1}, s_n = s' | s_0 = s) \quad (7)$$

$$= P(s_{n+1}|s_0 = s) \quad (8)$$

## 2 Markov Reward Process

**Definition 2.1** (Markov Reward Process). A Markov Reward Process (MRP) is a tuple  $(S, P, R, \gamma)$ , where

- $S$  is the finite state-space of the Markov process, (assume  $|S| < \infty$ )
- $P$  is the state transition probability model where  $P_{ss'} = P[s_{t+1} = s' | s_t = s]$
- $R : S \mapsto \mathbb{R}$  is a reward function that maps states to rewards,  $R(s) = E[r_t | s_t = s]$
- $\gamma \in [0, 1]$  is a discount factor

In a Markov reward process, whenever a transition happens from a current state  $s$  to a successor state  $s'$ , a reward is obtained depending on the current state  $s$ . Thus for the Markov process  $(s_0, s_1, \dots)$ , each transition  $s_t \rightarrow s_{t+1}$  is accompanied by a reward  $r_t$  for all  $i = 0, 1, \dots$ , and so a particular episode of the Markov reward process is represented as  $(s_0, r_0, s_1, r_1, s_2, r_2, \dots)$ . We should note that these rewards can be either deterministic or stochastic.

**Definition 2.2** (Expected reward). For a state  $s \in S$ , we define the expected reward  $R(s)$  by

$$R(s) = E[r_0 | s = s_0] \quad (9)$$

Just like the assumption of stationary transition probabilities, going forward we will also assume *stationarity of the rewards*. In the deterministic case, this implies that  $r_i = r_j$  wherever  $s_i = s_j$ . In the stochastic case, we require that the cumulative distribution functions (CDF) of the rewards conditioned on the current state be time independent:

$$F(r_i | s_i = s) = F(r_j | s_j = s), \quad (10)$$

where  $F$  denotes the cumulative distribution function of  $r_i$  conditioned on  $s_i$ . Therefore, the reward function  $R(s)$  is independent of  $t$ ,

$$R(s) = E(r_t | s_t = s) \quad (11)$$

**Definition 2.3** (Horizon). The horizon  $H$  of a Markov reward process is defined as the number of time steps in each episode (realization) of the process. The horizon can be finite or infinite. If the horizon is finite, then the process is also called a finite Markov reward process.

**Definition 2.4** (Return). The return  $G_t$  of a Markov reward process is defined as the discounted sum of rewards starting at time  $t$  up to the horizon  $H$ , and is given by

$$G_t = \sum_{k=t}^{H-1} \gamma^{k-t} r_k, \text{ for } t \in [0, H-1] \quad (12)$$

For example,  $G_0 = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots + \gamma^{H-1} r_{H-1}$

**Definition 2.5** (State value function). The state value function  $V_t(s)$  for a Markov reward process and a state  $s \in S$  is defined as the expected return starting from state  $s$  at time  $t$ , and is given by the following expression:

$$V_t(s) = E[G_t | s_t = s] \quad (13)$$

**Lemma 2.1.** Let us assume that

- Transition probability is stationary
- Rewards are stationary
- $H$  is infinite.

Then  $V_t(s)$  is independent of  $t$ . That is,

$$V_i(s) = V_j(s), \forall i, j \quad (14)$$

*Proof.* Even though this property seems obvious and intuitive, the proof is not totally straightforward (at least to me). I conduct the demonstration in two steps:

(1) I prove that  $E[r_{t+n} | s_t = s] = E[r_n | s_0 = s]$  by recursion on  $n$ :

- For  $n = 0$ , I use the result proved in lemma 1.1 to show that

$$E(r_t | s_t = s) = \sum_{r \in \mathbb{R}} r p(r_t = r | s_t = s) \quad (15)$$

$$= \sum_{r \in \mathbb{R}} r p(r_0 = r | s_0 = s) \quad (16)$$

$$= E(r_0 | s_0 = s) \quad (17)$$

- I assume  $E[r_{t+n} | s_t = s] = E[r_n | s_0 = s]$  for all  $n$ . Then, I show that  $E[r_{t+n+1} | s_t = s] = E[r_{n+1} | s_0 = s]$ .

$$E[r_{t+n+1} | s_t = s] = E[E[r_{t+n+1} | s_t = s, s_{t+n+1} = s'] | s_t = s] \quad (18)$$

$$= E[E[r_{t+n+1} | s_{t+n+1} = s'] | s_t = s] \quad (19)$$

$$= \sum_{s'} E[r_{t+n+1} | s_{t+n+1} = s'] P(s_{t+n+1} = s' | s_t = s) \quad (20)$$

$$= \sum_{s'} E[r_{t+n+1} | s_{t+n+1} = s'] P(s_{n+1} = s' | s_0 = s) \quad (21)$$

$$= \sum_{s'} E[r_{n+1} | s_{n+1} = s'] P(s_{n+1} = s' | s_0 = s) \quad (22)$$

$$= E[r_{n+1} | s_0 = s] \quad (23)$$

However, I think equation 18 does not come from the law of iterated expectation. That is what was confusing me (I saw this written everywhere). In fact, I think it comes from the tower property for two random variables, which states that

$$E[E[X|Y, Z] | Y] = E[X|Y], \quad (24)$$

whereas the law of iterated expectation is

$$E[E[X|Z]] = E[X]. \quad (25)$$

(2) Finally, I conclude that  $V_{t+n}(s) = V_t(s)$

$$V_{t+n}(s) = E[G_{t+n} | s_{t+n} = s] \quad (26)$$

$$= E\left[\sum_{k=0}^{\infty} \gamma^k r_{t+n+k} \mid s_{t+n} = s\right] \quad (27)$$

$$= \sum_{k=0}^{\infty} \gamma^k E[r_{t+n+k} \mid s_{t+n} = s] \quad (28)$$

$$= \sum_{k=0}^{\infty} \gamma^k E[r_{t+k} \mid s_t = s] \quad (29)$$

$$= V_t(s) \quad (30)$$

### 3 Other useful tricks

if  $X \geq 0$

$$E[X] = \int_0^{\infty} 1 - F_X(x) dx \quad (31)$$

### 4 Acknowledgments

I would like to specifically thank former CS 234 Teacher Assistant Rahul Sarkar for providing his quality notes that greatly helped me understand the material, and also inspired me to write my own notes. I used his work as a reference, and most of the changes made focused on deriving detailed proofs that were not always provided in his notes.