

**RL notes**

## MARKOV DECISION PROCESSES (MDP)

The at time  $t$  the state is  $S_t \in \mathcal{S}$ , the action is  $A_t \in \mathcal{A}$ , and the reward (received before seeing the state and doing the action) is  $R_t \in \mathcal{R}$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space, and  $\mathcal{R} \subset \mathbb{R}$  is the reward space.

The **trajectory** is

$$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \dots \quad (1)$$

The probability of getting to state  $s'$  and getting reward  $r$  after taking action  $a$  in state  $s$  is well defined, and given by

$$p(s', r \mid s, a) = \Pr \{S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a\} \quad (2)$$

where the function  $p$  defines **dynamics** of the MDP, and is called the **dynamics function**. The probabilities given  $p$  completely characterize the environment's dynamics. This also confirms that an MDP has the Markov property.

We can obtain the **state-transition probabilities** with

$$p(s' \mid s, a) = \Pr(S_t = s' \mid S_{t-1} = s, A_{t-1} = a) = \sum_{r \in \mathcal{R}} p(s', r \mid s, a) \quad (3)$$

And the expected rewards for a state-action pair:

$$r(s, a) = \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r \mid s, a) \quad (4)$$

and the expected reward given the next state:

$$r(s, a, s') = \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a, S_t = s'] = \sum_{r \in \mathcal{R}} r \frac{p(s', r \mid s, a)}{p(s' \mid s, a)} \quad (5)$$

## RETURNS AND EPISODES

The **return** is the discounted sum of rewards (if we are using discounting).

In an episodic task, this is

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T \quad (6)$$

where  $T$  is the final time step. We think of each episode ending in the **same** terminal state - can be thought of as an artificial state that occurs right *after* the real terminal state of the episode.

The final reward is given in this final terminal state.

For continuing tasks, the return is

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (7)$$

If we define the reward to be zero after the final state, this also holds for episodic tasks.

There is a recursive relationship between  $G_t$  and  $G_{t+1}$ :

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (8)$$

$$= R_t + \sum_{k=1}^{\infty} \gamma^k R_{t+k+1} \quad (9)$$

$$= R_t + \gamma \sum_{k=1}^{\infty} \gamma^{k-1} R_{t+k+1} \quad (10)$$

$$= R_t + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} \quad (11)$$

$$= R_t + \gamma G_{t+1} \quad (12)$$

$$(13)$$

## POLICIES AND VALUE FUNCTION

The **value** function is defined to be

$$v_{\pi}(s) = \mathbb{E}_{\pi} [G_t \mid S_t = s] \quad (14)$$

$$= \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right] \quad (15)$$

and the **action-value** function is

$$q_{\pi}(s, a) = \mathbb{E}_{\pi} [G_t \mid S_t = s, A_t = a] \quad (16)$$

$$= \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right] \quad (17)$$

$$(18)$$

The **Bellman equation** for the value function is

$$v_{\pi}(s) = \mathbb{E}_{\pi} [G_t \mid S_t = s] \quad (19)$$

$$= \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} \mid S_t = s] \quad (20)$$

$$= \sum_a \pi(a \mid s) \sum_{s'} \sum_r p(s', r \mid s, a) [r + \gamma \mathbb{E}_{\pi} [G_{t+1} \mid S_{t+1} = s']] \quad (21)$$

$$= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_{\pi}(s')] \quad (22)$$

$$(23)$$