

Reinforcement Learning Equations

Aaron Hao Tan

Finite Markov Decision Processes

Components of MDP: Decision epochs, State, Action, Transition probability, Rewards

$$\{T, S, A_s, p_t(\cdot|s, a), r_t(s, a)\} \quad (1)$$

A state is *Markov* if and only if

$$\mathbb{P}[S_{t+1}|S_t] = \mathbb{P}[S_{t+1}|S_1, \dots, S_t] \quad (2)$$

State-transition probabilities (p characterize the environment's dynamics). More often used than equations 4 and 5.

$$\begin{aligned} p(s'|s, a) &= Pr[S_t = s'|S_{t-1} = s, A_{t-1} = a] \\ &= \sum_{r \in \mathcal{R}} p(s', r|s, a) \end{aligned} \quad (3)$$

Expected rewards for state-action pairs and for state-action-next-state triples. The current reward depends on the previous state and action.

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

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

Policy

$$\pi(a|s) = Pr(A_t = a|S_t = s) \quad (6)$$

Returns (sum of rewards)

$$G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T \quad (7)$$

Discounted Returns (γ = discount rate). If $\gamma < 1$, the infinite sum in equation 8 would have a finite value. If $\gamma = 0$, the agent is concerned with only maximizing immediate rewards (myopic). If the reward is +1, the return is equivalent to $G_t = \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}$.

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (8)$$

$$\begin{aligned} G_t &\doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned} \quad (9)$$

State Value Function. Value of terminal state (if any) is always zero.

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t|S_t = s] = \mathbb{E}_{\pi}[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}|S_t = s] \quad \forall s \in S \quad (10)$$

Bellman Equation for v_{π} : the value of the start state must equal the (discounted) value of the expected next state, plus the reward expected along the way.

$$\begin{aligned} v_{\pi}(s) &= E_{\pi}[R_{t+1} + \gamma \cdot G_{t+1}|S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) [r + \gamma \mathbb{E}_{\pi}[G_{t+1}|S_{t+1} = s']] \\ &= \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma v_{\pi}(s')] \quad \forall s \in S \end{aligned} \quad (11)$$

Action Value Function

$$q_\pi(s, a) \doteq \mathbb{E}_\pi [G_t | S_t = s, A_t = a] = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right] \quad (12)$$

Bellman Equation for $q_\pi(s, a)$

$$q_\pi(s, a) = \sum_{s', r} p(s', r | s, a) \left[r + \gamma \sum_{a'} \pi(a' | s') q_\pi(s', a') \right] \quad (13)$$

The following figure represents the state value and action state value functions graphically.

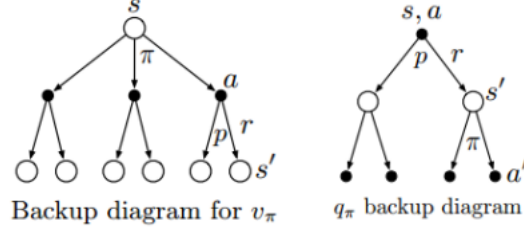


Figure 1: Value and State-Value Functions

Optimal State Value Function. A policy π is better or equal to another policy π' if its expected return is greater than or equal to that of π' for all states ($\pi \geq \pi'$ if and only if $v_\pi(s) \geq v_{\pi'}(s)$).

$$v_*(s) \doteq \max_{\pi} v_\pi(s) \quad (14)$$

Optimal Action Value Function and its relation to v_* . Equation 15 gives the expected return for taking action a in state s and thereafter following an optimal policy. Comparing with equation 11, $v_\pi(s)$ represents the value of a state as the summation of future returns. In equation 15, $v_*(S_{t+1})$ gives the optimal value of all future states from the given state and action.

$$\begin{aligned} q_*(s, a) &\doteq \max_{\pi} q_\pi(s, a) \\ &= \mathbb{E} [R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \end{aligned} \quad (15)$$

Relationship between q_* and v_* . The value of a state under an optimal policy must equal the expected return for the best action from that state. Second last line of equation 16 shows that $G_t = R_{t+1} + \gamma v_*(S_{t+1})$ when following π_* . The last line is the same as equation 17. For finite MDPs, the last line of equation 16 has a unique solution.

$$\begin{aligned} v_*(s) &= \max_{a \in A(s)} q_*(s, a) \\ &= \max_a q_*(s, a) \\ &= \max_a \mathbb{E}_{\pi_*} [G_t | S_t = s, A_t = a] \\ &= \max_a \mathbb{E} [R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')] \end{aligned} \quad (16)$$

Bellman Optimality Equations

$$\begin{aligned} v_*(s) &= \max_a \mathbb{E} [R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')] \end{aligned} \quad (17)$$

$$\begin{aligned} q_*(s, a) &= \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') | S_t = s, A_t = a \right] \\ &= \sum_{s', r} p(s', r | s, a) \left[r + \gamma \max_{a'} q_*(s', a') \right] \end{aligned} \quad (18)$$

In v_* , the next best action is selected based on the expected reward and value of future states. In q_* , the state and action is given as well as the reward. From the new state s' , the best action is chosen.

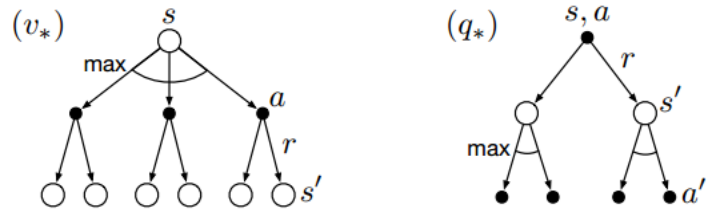


Figure 2: Bellman Optimality Equations