

---

# Reinforcement Learning

---

Zhang Fengchen

## Abstract

Leave it blank for now.

## 1. Background

### (1) Probability Definition

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

$$\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 1 \quad \forall s \in S, a \in A$$

$$p(s' | s, a) = \sum_{r \in R} p(s', r | s, a)$$

$$r(s, a) = \sum_{s' \in S} \sum_{r \in R} (r * p(s', r | s, a))$$

$$r(s, a, s') = \frac{\sum_{r \in R} (r * p(s', r | s, a))}{p(s' | s, a)}$$

### (2) Bellman Equations

$v(s)$ :

$$v_\pi(s) = E_\pi [G_t | S_t = s] \quad \forall s \in S$$

$$v_\pi(s) = E_\pi [R_{t+1} + \gamma G_{t+1} | S_t = s] \quad \forall s \in S$$

$$v_\pi(s) = \sum_{a \in A} (\pi(a | s) * q_\pi(s, a)) \quad \forall s \in S$$

$$v_\pi(s) = \sum_{a \in A} \pi(a | s) \sum_{s', r} p(s', r | s, a) * [r + \gamma E_\pi [G_{t+1} | S_{t+1} = s']] \quad \forall s \in S$$

$$v_\pi(s) = \sum_{a \in A} \pi(a | s) \sum_{s', r} p(s', r | s, a) * [r + \gamma v_\pi(s')] \quad \forall s \in S$$

$q(s, a)$ :

$$q_\pi(s, a) = E_\pi [G_t | S_t = s, A_t = a] \quad \forall s \in S, a \in A$$

$$q_\pi(s, a) = E_\pi [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \quad \forall s \in S, a \in A$$

$$q_\pi(s, a) = \sum_{s', r} p(s', r | s, a) * [r + \gamma E_\pi [G_{t+1} | S_{t+1} = s']]$$

$$q_\pi(s, a) = \sum_{s', r} p(s', r | s, a) * [r + \gamma v_\pi(s')] \quad \forall s \in S, a \in A$$

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

Optimal Equations:

$$v_*(s) = \max_{a \in A} q_{\pi_*}(s, a) \quad \forall s \in S$$

### (3) Dynamic Programming

Policy Improvement Theorem:

$$E_{\pi'} [q_\pi(s, \pi'(s))] \geq v_\pi(s) \quad \forall s \in S$$

$$\pi' \geq \pi \quad \longleftrightarrow \quad v_{\pi'}(s) \geq v_\pi(s) \quad \forall s \in S$$

Policy Evaluation:

$$v_{k+1}(s) = E_\pi [R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s]$$

$$= \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_k(s')]$$

Policy Improvement:

$$\pi'(s) = \arg \max_a q_\pi(s, a)$$

Value Iteration:

$$v_{k+1}(s) = \max_a \mathbb{E} [R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s, A_t = a]$$

## 2. Proposed Solution

### (1) Frozen Lake

## 3. Numerical Results

Your experiment results should be here. You should add the figure/table if necessary.

---

. Correspondence to: Zhang Fengchen  
<fczhang0606@gmail.com>.