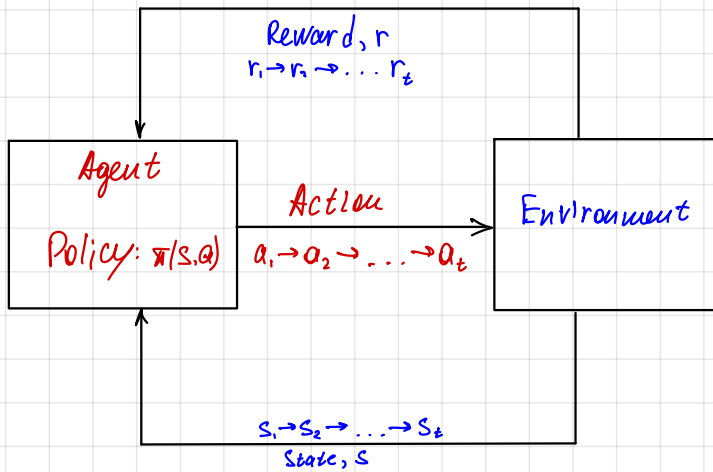


1. Overview

Policy: $\pi(s, a) = \Pr(a=a | s=s)$



Value: $V_\pi(s) = \mathbb{E}(\sum_t \gamma^t r_t | s_0 = s)$

2. Q-learning

$Q(s, a)$ = Quality of state / action pair

- combine Value $V_\pi(s)$ and Policy $\pi(s, a)$.

$$Q^{\text{update}}(s_t, a_t) = Q^{\text{old}}(s_t, a_t) + \alpha (r_t + \gamma \max_a Q(s_{t+1}, a) - Q^{\text{old}}(s_t, a_t))$$