

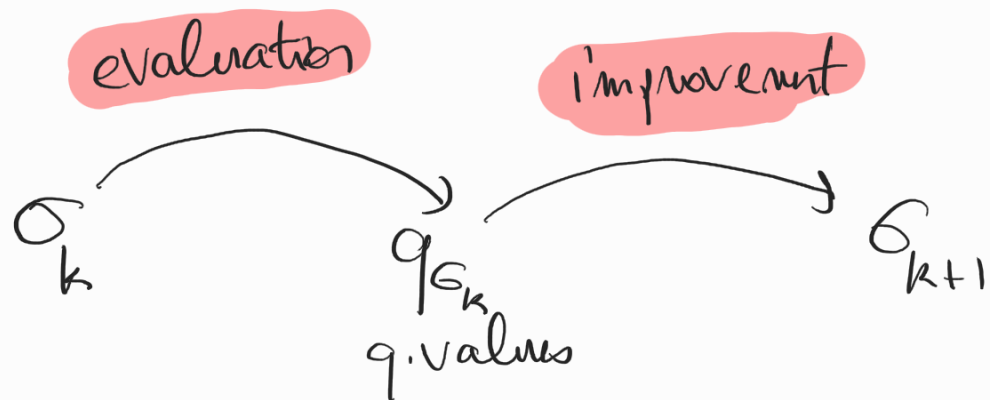
Value iteration $v: S \rightarrow \mathbb{R}$

$v_0 \rightarrow v_1 \rightarrow v_2 \rightarrow v_3 \dots$

policy iteration
||
strategy improvement

$\pi_0 \rightarrow \pi_1 \rightarrow \pi_2 \rightarrow \pi_3 \dots$

$\pi: S \rightarrow A$



improvement:

$$\pi_{k+1}(s) = \operatorname{argmax}_{a \in A} q_{\pi_k}(s, a)$$

evaluation:

$$\pi \longrightarrow q_\pi$$

On-policy Bellman equations for q -values:

$$q_\pi(s, a) = \sum_{s', r} \Delta(s, a)(s', r) \left(r + \gamma q_\pi(s', \pi(s')) \right)$$

$s' \in S$
 $r \in \mathbb{R}$

$val_{\pi}(s')$

"policy evaluation" : (variant of VI)
where π is fixed

$$q_{\pi}^0(s, a) = 0 \quad \text{for all } s, a$$

For $k = 0, \dots$:

$$q_{\pi}^{k+1}(s, a) = \sum_{\substack{s' \in S \\ r \in \mathbb{R}}} \Delta(s, a)(s', r) \left(r + \gamma q_{\pi}^k(s', \pi(s')) \right)$$

until: $\|q_{\pi}^{k+1} - q_{\pi}^k\| \leq \epsilon$

PI:

π_0 arbitrary policy

iterate:

• evaluation (π_k) $\rightarrow q_{\pi_k}$

• improve (q_{π_k})

• $b_{k+1} = \text{improve}(b_k, \gamma b_k)$

MDPs

Bellman equations

↳ Value iteration

↳ Policy iteration