

Policy Optimization

Suppose we aim to find an optimal policy π^* , i.e.,

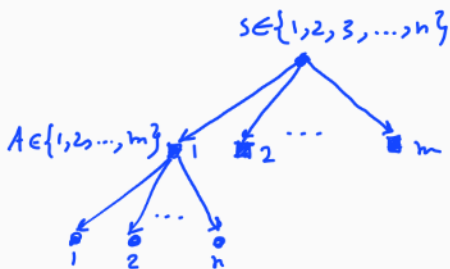
$$\begin{aligned} & V^{\pi^*}(s) \geq V^{\pi}(s) \quad \forall s \in S, \forall \pi \in \Pi \\ \text{or} \quad & V^{\pi^*} \geq V^{\pi} \quad \forall \pi \in \Pi \\ \text{or} \quad & V^{\pi^*} = V^* = \max_{\pi \in \Pi} V^{\pi} \end{aligned}$$

Bellman Optimality Equations:

The optimal value functions, i.e., V^* and Q^* , satisfy the following equations:

$$V^*(s) = \max_{a \in A} \left[R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^*(s')] \right] \quad \forall s \in S$$

$$Q^*(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\max_{a' \in A} Q^*(s', a') \right] \quad \forall s \in S, \forall a \in A$$



$$V^{\pi}(s) = \mathbb{E}_{a \sim \pi} [R(s, a) + \gamma \mathbb{E}_{s'} [V^{\pi}(s')]] \quad \text{Bellman consistency equation}$$

$$V^* = \max_{a \in A} [R(s, a) + \gamma \mathbb{E}_{s'} [V^*(s')]] \quad \text{Bellman optimality equation}$$

Theorem [Bellman optimality]: π^* is an optimal policy if and only if v^{π^*} satisfies the Bellman optimality equation.

Proof (sketch):

① π^* an optimal policy $\Rightarrow v^{\pi^*}$ satisfies the Bellman optimality equation

Consider the following deterministic, stationary policy

$$\tilde{\pi}(s) = \arg \max_{a \in A} [R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [v^{\pi^*}(s')]] \quad \forall s \in S$$

One can show that $v^{\pi^*} = v^{\tilde{\pi}}$ and $v^{\tilde{\pi}}$ satisfies the Bellman optimality equation.

② v^{π} satisfies the Bellman optimality equation $\Rightarrow \pi$ an optimal policy

Consider an optimal policy π^* .

By ①, we know that v^{π^*} satisfies the Bellman optimality equation.

One can show that $|v^{\pi}(s) - v^{\pi^*}(s)| \leq 0$ for all $s \in S$.

Corollary: The deterministic stationary policy

$$\pi(s) = \operatorname{argmax}_{a \in A} \left[\underbrace{R(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^*(s')]}_{Q^*(s, a)} \right]$$

$$\text{or } \pi(s) = \operatorname{argmax}_{a \in A} Q^*(s, a)$$

is an optimal policy.

The optimal value function satisfies the Bellman optimality equation:

$$V^*(s) = \max_{a \in A} \left[R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^*(s')] \right], \quad \forall s \in S$$

Value Iteration

VI aims to find (approximate) the optimal value function v^* through the fixed-point iteration algorithm.

Value Iteration Algorithm

- Initialize \vec{v}_0
- For $t = 0, 1, 2, \dots, T-1$:

$$v_{t+1}(s) = \max_{a \in A} [R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [v_t(s')]] \quad \forall s \in S$$

- Return $\vec{v}_T \approx \vec{v}^*$

or

$$\pi_T(s) = \arg \max_{a \in A} [R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [v_T(s')]] \quad \forall s \in S$$

$v^{\pi_T} \approx v^*$

Bellman optimality operator: $\mathcal{T} : \mathbb{R}^{|S|} \rightarrow \mathbb{R}^{|S|}$

$$(\mathcal{T}v)(s) = \max_{a \in A} [R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [v(s')]] \quad \forall s \in S$$