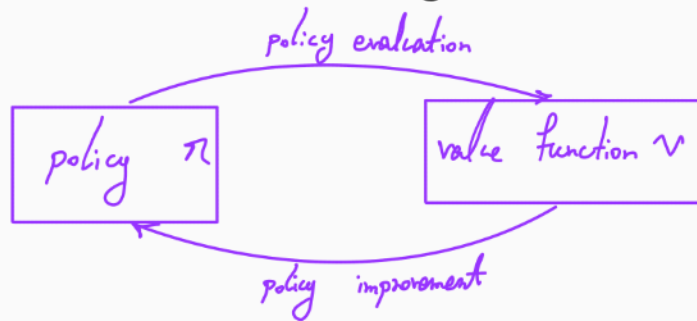


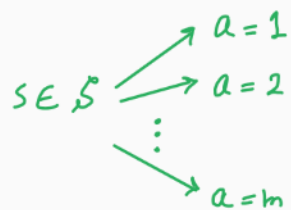
Policy Iteration

PI aims to find an (approximately) optimal policy π^* through iterative policy evaluation and policy improvement.



Policy Improvement: Suppose there exists a deterministic, stationary policy π with value function v^π .

***Question:** Would changing the policy at a state improve the policy?



$$\pi(s) = a$$

Would setting $\pi(s) = a' \neq a$ be helpful?

Recall

$$Q^\pi(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [v^\pi(s')]$$

If $\exists a \in A$ s.t. $Q^\pi(s, a) > v^\pi(s)$, then changing the policy at state s to a improves the policy.

Note: If for all $s \in S$, $\max_{a \in A} Q^\pi(s, a) = v^\pi(s)$, then

$$v^\pi(s) = v^*(s) \quad \forall s \in S \quad \pi \text{ an optimal policy}$$

Policy Improvement Theorem: Let π and π' be a pair of deterministic, stationary policies such that

$$Q^\pi(s, \pi'(s)) \geq v^\pi(s) \quad \forall s \in S.$$

Then, π' must be as good as or better than π , i.e.,

$$v^{\pi'} \geq v^\pi.$$

Policy Iteration Algorithm

- Initialize $\pi_0: S \rightarrow A$
- For $t = 0, 1, 2, \dots, T-1$:

policy evaluation . Evaluate π_t by computing v^{π_t}

policy improvement . Improve the policy by

$$\pi_{t+1}(s) = \underset{a \in A}{\operatorname{argmax}} \left[R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [v^{\pi_t}(s')] \right] \quad \forall s \in S$$

- Return π_T

Lemma: Let v^{π_t} represent the value function of policy π_t at iteration t of PI. It holds that

$$v^{\pi_{t+1}} \geq v^{\pi_t}$$

element-wise \leftarrow

Monotonic Improvement

proof :

$$\begin{aligned}
 \forall s \in \mathcal{S} : \quad & v^{\pi_{t+1}}(s) - v^{\pi_t}(s) = R(s, \pi_{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi_{t+1}(s))} [v^{\pi_{t+1}}(s')] \\
 & - (R(s, \pi_t(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi_t(s))} [v^{\pi_t}(s')]) \\
 & \geq R(s, \pi_{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi_{t+1}(s))} [v^{\pi_{t+1}}(s')] \\
 & - (R(s, \pi_t(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi_t(s))} [v^{\pi_t}(s')]) \\
 & = \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi_{t+1}(s))} [v^{\pi_{t+1}}(s') - v^{\pi_t}(s')] \\
 \longrightarrow & v^{\pi_{t+1}} - v^{\pi_t} \geq \gamma P^{\pi_{t+1}} (v^{\pi_{t+1}} - v^{\pi_t}) \\
 & \geq \gamma P^{\pi_{t+1}} (\gamma P^{\pi_{t+1}} (v^{\pi_{t+1}} - v^{\pi_t})) = \gamma^2 (P^{\pi_{t+1}})^2 (v^{\pi_{t+1}} - v^{\pi_t}) \\
 & \vdots \\
 & \geq \gamma^k (P^{\pi_{t+1}})^k (v^{\pi_{t+1}} - v^{\pi_t})
 \end{aligned}$$

$$\text{Let } k \longrightarrow \infty \implies v^{\pi_{t+1}} - v^{\pi_t} \geq 0 \implies v^{\pi_{t+1}} \geq v^{\pi_t}$$

■

Theorem: The value of the final policy v^{π_T} returned by the PI algorithm satisfies $\|v^{\pi_T} - v^*\|_\infty \leq \gamma^T \|v^{\pi_0} - v^*\|_\infty$.

$$\begin{aligned}
\text{Proof: } v^*(s) - v^{\pi_t}(s) &= \max_{a \in A} \left[R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [v^*(s')] \right] \\
&\quad - \left(R(s, \pi_t(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi_t(s))} [v^{\pi_t}(s')] \right) \\
&\leq \max_{a \in A} \left[R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [v^*(s')] \right] \\
&\quad - \left(R(s, \pi_t(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi_t(s))} [v^{\pi_{t-1}}(s')] \right) \\
&= \max_{a \in A} \left[R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [v^*(s')] \right] \\
&\quad - \max_{a \in A} \left[R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [v^{\pi_{t-1}}(s')] \right]
\end{aligned}$$

$$\begin{aligned}
\Rightarrow \forall s \in S: |v^*(s) - v^{\pi_t}(s)| &\leq \left| \max_{a \in A} \left[R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [v^*(s')] \right] \right. \\
&\quad \left. - \max_{a \in A} \left[R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [v^{\pi_{t-1}}(s')] \right] \right| \\
&\leq \max_{a \in A} \left| \left(R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [v^*(s')] \right) \right. \\
&\quad \left. - \left(R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [v^{\pi_{t-1}}(s')] \right) \right| \\
&= \gamma \max_{a \in A} \left| \mathbb{E}_{s' \sim P(\cdot | s, a)} [v^*(s') - v^{\pi_{t-1}}(s')] \right| \\
&\leq \gamma \max_{a \in A} \mathbb{E}_{s' \sim P(\cdot | s, a)} [|v^*(s') - v^{\pi_{t-1}}(s')|] \\
&\leq \gamma \max_{s' \in S} |v^*(s') - v^{\pi_{t-1}}(s')|
\end{aligned}$$

$$\Rightarrow \max_s |v^*(s) - v^{\pi_t}(s)| \leq \gamma \max_{s \in S} |v^*(s) - v^{\pi_{t-1}}(s)|$$

$$\begin{aligned} \Rightarrow \|v^{\pi_t} - v^*\|_\infty &\leq \gamma \|v^{\pi_{t-1}} - v^*\|_\infty \\ &\leq \gamma^2 \|v^{\pi_{t-2}} - v^*\|_\infty \\ &\vdots \\ &\leq \gamma^t \|v^{\pi_0} - v^*\|_\infty \end{aligned}$$

Value Iteration

vs.

Policy Iteration

- Initialize v_0
- Each iteration:
Apply Bellman optimality operator
- Return v_T & π_T

- Initialize π_0
- Each iteration:
Apply policy evaluation
Apply policy improvement
- Return π_T

search over the space of

value functions, $v \in \mathbb{R}^{|S|}$

deterministic stationary policies, $|A|^{|S|}$ policies

convergence rate of $\|v^{\pi_T} - v^*\|_\infty$

geometric in T

geometric in T

computational complexity

polynomial, $O\left(\frac{|S|^2 |A|}{1-\gamma} \log\left(\frac{1}{(1-\gamma)\epsilon}\right)\right)$

strongly polynomial, $O\left(\frac{|S|^4 |A|^2}{1-\gamma} \log\left(\frac{1}{1-\gamma}\right)\right)$

Finite-Horizon MDP

Goal: The learning agent aims to find a policy π that maximizes the expected cumulative reward over a finite horizon

$$\pi^* \in \arg \max_{\pi \in \Pi} \mathbb{E}_{\substack{s_0 \sim \mathcal{M}_0 \\ A_t \sim \pi(z_t) \\ s_{t+1} \sim P(\cdot | s_t, A_t)}} \left[\overbrace{\sum_{t=0}^{T-1} R(s_t, a_t)}^{\text{return}} \right]$$

γ^t

Value functions

- (state) value function of a policy π at time $t \in \{0, 1, \dots, T-1\}$,

$V_t^\pi: S \rightarrow \mathbb{R}$, is defined as

$$V_t^\pi(s) = \mathbb{E}_{\substack{A_t \sim \pi(z_t) \\ s_{t+1} \sim P(\cdot | s_t, A_t)}} \left[\sum_{t'=t}^{T-1} R(s_{t'}, A_{t'}) \mid s_t = s \right] \quad \forall s \in S$$

$\gamma^{t'-t}$

- (state-) action value function of a policy π at time $t \in \{0, 1, \dots, T-1\}$,

$Q_t^\pi: S \times A \rightarrow \mathbb{R}$, is defined as

$$Q_t^\pi(s, a) = \mathbb{E}_{\substack{A_t \sim \pi(z_t) \\ s_{t+1} \sim P(\cdot | s_t, A_t)}} \left[\sum_{t'=t}^{T-1} R(s_{t'}, A_{t'}) \mid s_t = s, A_t = a \right] \quad \forall s \in S, \forall a \in A$$

$\gamma^{t'-t}$

Theorem: There exists an optimal, deterministic, history-independent, time-varying policy π^* , i.e.,

$$\exists \pi^* = (\pi_0^*, \pi_1^*, \dots, \pi_{T-1}^*), \quad \pi_t^* : S \rightarrow A$$

$$\text{s.t.} \quad V_t^{\pi^*}(s) \geq V_t^\pi(s) \quad \forall s \in S, \forall \pi \in \Pi$$

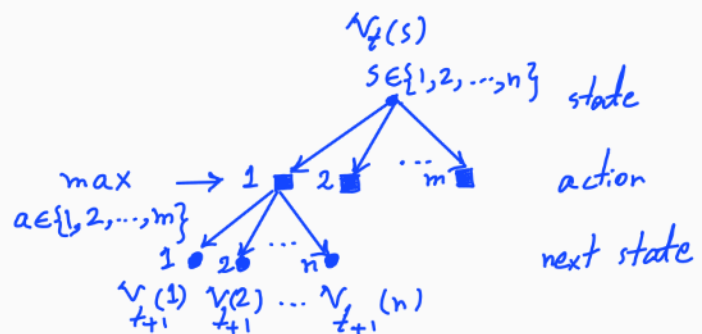
→ we look for finding an optimal policy $\pi^* = (\pi_0^*, \pi_1^*, \dots, \pi_{T-1}^*)$, which is a sequence of deterministic policies dependent on state and time.

* Bellman (consistency) Equations:

$$V_t^\pi(s) = \mathbb{E}_{a \sim \pi_t(s)} \left[R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V_{t+1}^\pi(s')] \right]$$

$$Q_t^\pi(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V_{t+1}^\pi(s')]$$

* Bellman Optimality Equations:



$$V_t^* = \max_{a \in A} \left[R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V_{t+1}^*(s')] \right]$$

$$Q_t^*(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\max_{a' \in A} Q_{t+1}^*(s', a') \right]$$

Note: The optimal value function is also a sequence, i.e.,

$$V^* = (V_0^*, V_1^*, \dots, V_T^*)$$

objective \swarrow \searrow 0

Value Iteration (for finite-Horizon MDP)

VI aims to find the (exact) optimal value function and an optimal policy through dynamic programming (backward induction).

Iterative method:

- Initialize $V_T^* = \mathbf{0}_{|S| \times 1}$

- For $t \in \{T-1, T-2, \dots, 0\}$:

$$V_t^*(s) = \max_{a \in A} [R_t(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V_{t+1}^*(s')]] \quad \forall s \in S$$

$$\pi_t^*(s) = \operatorname{argmax}_{a \in A} [R_t(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V_{t+1}^*(s')]] \quad \forall s \in S$$

- Return $V^* = (V_0^*, V_1^*, \dots, V_T^*)$ if time-varying

and $\pi^* = (\pi_0^*, \pi_1^*, \dots, \pi_{T-1}^*)$

Note: These expressions and the VI solution for finite-horizon MDP can be extended to the case of time-varying MDPs where the transition P_t and reward R_t are time-dependent.

Effective Horizon

Question: Can we approximate an infinite-horizon discounted return with a finite-horizon (truncated) return?

$$G = \underbrace{\sum_{t=0}^{\infty} \gamma^t R(s_t, A_t)}_{\text{true return } G} = \underbrace{\sum_{t=0}^{T-1} \gamma^t R(s_t, A_t)}_{\text{truncated return } \hat{G}} + \underbrace{\sum_{t=T}^{\infty} \gamma^t R(s_t, A_t)}_{\text{error } G - \hat{G}}$$

$$\text{If } |R(s, a)| \leq R_{\max} \quad \forall s \in \mathcal{S}, \forall a \in \mathcal{A}$$

$$\leq \gamma^T \sum_{t=T}^{\infty} \gamma^{t-T} R_{\max}$$

$$\longrightarrow |\text{error}| = |G - \hat{G}| \leq \frac{\gamma^T R_{\max}}{1 - \gamma} \leq \epsilon$$

$$\longrightarrow T \geq \frac{\log \frac{R_{\max}}{\epsilon(1-\gamma)}}{\log \frac{1}{\gamma}} \leftarrow \geq 1-\gamma$$

$$\xrightarrow[\text{simplicity}]{\text{for}} T_{\gamma, \epsilon} := \frac{\log \frac{R_{\max}}{\epsilon(1-\gamma)}}{1 - \gamma}$$

effective horizon

Note: The notion of effective horizon can be used to create approximations between infinite-horizon discounted MDPs and finite-horizon MDPs.