

In a realistic setting, an agent in a RL framework does not have a priori knowledge of the environment.

We can consider a simplified setting, however, in which both the rules of transition and reward system are known to the agent.

More precisely, we may suppose that the following are known:

(a.) $P(s' | s, a)$, i.e. the probability of transitioning to state s' , given that the agent's present state is $s \in S$ and action $a \in A$ is taken.

(b.) $E(r(s, a))$, i.e. the expected reward from taking action a in state s .

Since we don't have to estimate these environmental parameters, the RL framework reduces to a planning problem in this setting.

We consider three methods for solving this planning problem:

(a.) Value iteration → Use the Bellman equations and iteratively determine $V^*(s), \forall s \in S$

(b.) Policy iteration → Use policy evaluation and our matrix inversion result from last lecture

(c.) Linear programming → Formulate the problem as maximizing a linear objective subject to linear constraints, then solve with, e.g., Simplex method.

Value Iteration

Let $V \in \mathbb{R}^{|S|}$ and let $\Phi = \mathbb{R}^{|S|} \rightarrow \mathbb{R}^{|S|}$ be defined as in the Bellman equations:

$$\forall s \in S, \Phi(V)(s) = \max_{a \in A} \left(E(r(s, a)) + \gamma \sum_{s' \in S} P(s' | s, a) \cdot V(s') \right)$$

This implicitly defines a policy, and we can thus re-write

$$\Phi(V) = \max_{\pi} (R_{\pi} + \gamma P_{\pi} V), \quad (\star)$$

where $P_{\pi} \in \mathbb{R}^{|S| \times |S|}$ is defined as $(P_{\pi})_{ss'} = P(s' | s, \pi(s))$, for all $s, s' \in S$, and where $R_{\pi} \in \mathbb{R}^{|S| \times 1}$ is defined as $(R_{\pi})_s = \mathbb{E}(r(s, \pi(s)))$, for all $s \in S$.

The main idea of value iteration is to iteratively apply Φ to an arbitrary starting point until some convergence criterion is achieved. That is, let $V_0 \in \mathbb{R}^{|S| \times 1}$ be arbitrary, and set $V_{i+1} = \Phi(V_i)$, $i \geq 0$, until the stopping condition $\|V_i - \Phi(V_i)\| < \frac{(1-\gamma)}{\gamma} \cdot \epsilon$ is reached for some $\epsilon > 0$.

This converges (though perhaps slowly) to the optimal value:

Theorem (convergence of value iteration): Let V^* be the value function associated to an optimal π^* . Then for any initial guess V_0 , the sequence $V_{i+1} = \Phi(V_i)$, $i \geq 0$ converges to V^* .

Proof: Note that it is sufficient to show Φ is contracting, i.e. has a Lipschitz constant < 1 w.r.t the $\|\cdot\|_{\infty}$ topology. In fact, we will show Φ is $\gamma < 1$ Lipschitz. Indeed, if we show this, then

$$\begin{aligned} & \|V^* - V_{n+1}\|_{\infty} \\ &= \|\Phi(V^*) - V_{n+1}\|_{\infty}, \quad \text{since } \Phi(V^*) = V^* \text{ by the Bellman conditions} \\ &= \|\Phi(V^*) - \Phi(V_n)\|_{\infty} \end{aligned}$$

$$\leq \gamma \|V^* - V_n\|_\infty$$

Inducting, we see $\|V^* - V_{n+1}\|_\infty \leq \gamma^{n+1} \cdot \|V^* - V_0\|_\infty$; taking $n \rightarrow \infty$ gives convergence.

It remains to show Φ is γ -Lipschitz. Well, for any $s \in S$, and any $V \in \mathbb{R}^{|S|}$, let $a^*(s)$ be the action maximizing $R_\pi + \gamma P_\pi V$, which defines $\Phi(V)(s)$. Then $\forall s \in S$ and $\forall U \in \mathbb{R}^{|S|}$, we have

$$\begin{aligned} \Phi(V)(s) - \Phi(U)(s) &\leq \Phi(V)(s) - (\mathbb{E}(r(s, a^*(s))) + \gamma \sum_{s' \in S} P(s' | s, a^*(s)) \cdot U(s')) \\ &= \gamma \sum_{s' \in S} P(s' | s, a^*(s)) \cdot (V(s') - U(s')) \\ &\leq \gamma \sum_{s' \in S} P(s' | s, a^*(s)) \cdot \|V - U\|_\infty \\ &= \gamma \cdot \|V - U\|_\infty \cdot \sum_{s' \in S} P(s' | s, a^*(s)) \\ &= \gamma \cdot \|V - U\|_\infty, \end{aligned}$$

as desired. \blacksquare (gets ϵ -close to optimality), and

Note that the proposed algorithm converges in $O(\log \frac{1}{\epsilon})$ iterations:

$$\begin{aligned} \|V^* - V_{n+1}\|_\infty &\leq \|V^* - \Phi(V_{n+1})\|_\infty + \|\Phi(V_{n+1}) - V_{n+1}\|_\infty \\ &= \|\Phi(V^*) - \Phi(V_{n+1})\|_\infty + \|\Phi(V_{n+1}) - \Phi(V_n)\|_\infty \\ &\stackrel{\gamma\text{-Lipschitz}}{\rightarrow} \leq \gamma \cdot \|V^* - V_{n+1}\|_\infty + \gamma \cdot \|V_{n+1} - V_n\|_\infty \end{aligned}$$

$$\Rightarrow (1 - \gamma) \cdot \|V^* - V_{n+1}\|_\infty \leq \gamma \cdot \|V_{n+1} - V_n\|_\infty$$

$$\begin{aligned} \Rightarrow \|V^* - V_{n+1}\|_\infty &\leq \frac{\gamma}{1-\gamma} \|V_{n+1} - V_n\|_\infty \\ &\leq \frac{\gamma}{1-\gamma} \cdot \frac{1-\gamma}{\gamma} \cdot \varepsilon \quad \text{by our stopping condition} \\ &\leq \varepsilon. \end{aligned}$$

(4)

To see convergence, note that if n is largest integer s.t. $\frac{(1-\gamma)\varepsilon}{\gamma} \leq \|V_{n+1} - V_n\|_\infty$, i.e. if the algorithm converges on the $(n+1)$ st step, then

$$\frac{(1-\gamma)\varepsilon}{\gamma} \leq \|V_{n+1} - V_n\|_\infty \leq \gamma^n \cdot \|\Phi(V_0) - V_0\|_\infty$$

$$\Rightarrow \frac{(1-\gamma)\varepsilon}{\gamma} \cdot \|\Phi(V_0) - V_0\|_\infty \leq \gamma^n$$

$$\Rightarrow \log\left(\frac{1-\gamma}{\gamma}\right) + \log(\varepsilon) + \log(\|\Phi(V_0) - V_0\|_\infty) \leq n \cdot \log(\gamma)$$

$$\Rightarrow \frac{O(\log(\frac{1}{\varepsilon}))}{\log(\gamma)} \geq n.$$

Policy Iteration

Instead of iterating on V , we can iterate on π . That is, we can pick an arbitrary initial $\pi = \pi_0$, then solve $(I - \gamma P_\pi)V = R_\pi$ to get V , then set $\pi = \arg\max_{\pi'} (R_{\pi'} + \gamma P_{\pi'} V)$. Iterating until some kind of criterion

for stopping is met yields an approximately optimal π .
 Indeed, let $\{V_n\}_{n=0}^\infty$ be the V sequence induced by the above scheme.

Thm (Convergence of Policy Iteration): Let $\{V_n\}_{n=1}^{\infty}$ be the sequence of policy values computed above. Then $\forall n, V_n \leq V_{n+1} \leq V^*$. (5)

Proof: Let π_{n+1} be the policy improvement at the n^{th} iteration of the algorithm. We claim that if $(Y-X) \geq 0$, for any $Y, X \in \mathbb{R}^{|\mathcal{S}|}$, then

$$(I - \gamma P_{\pi_{n+1}})^{-1} (Y-X) \geq 0 \quad (\star)$$

also. Well, we know from $P_{\pi_{n+1}}$ Markov that $\|\gamma P_{\pi_{n+1}}\|_{\infty} < 1$, so that we

may expand in Neumann series:

$$(I - \gamma P_{\pi_{n+1}})^{-1} = \sum_{k=0}^{\infty} (\gamma P_{\pi_{n+1}})^k$$

Noting that $P_{\pi_{n+1}}$ is pointwise non-negative yields (\star) .

Now, by definition of π_{n+1} , we have

$$\underbrace{R_{\pi_{n+1}} + \gamma P_{\pi_{n+1}} V_n}_{\text{by optimality of } \pi_{n+1}} \geq \underbrace{R_{\pi_n} + \gamma P_{\pi_n} V_n}_{\text{by construction}} = V_n$$

In particular,

$$R_{\pi_{n+1}} \geq (I - \gamma P_{\pi_{n+1}}) V_n$$

$$\Rightarrow (I - \gamma P_{\pi_{n+1}})^{-1} R_{\pi_{n+1}} \geq V_n$$

$$\underbrace{\hspace{10em}}_{\parallel}$$

$$V_{n+1}$$

Noting that V^* is maximal by definition gives the result.

If we stop once $\pi_n = \pi_{n+1}$, we may have to iterate over all $|A|^{|\mathcal{S}|}$ policies. In practice, bounds of the form $O\left(\frac{|A|^{|\mathcal{S}|}}{|\mathcal{S}|}\right)$ are possible. (6)

Policy & value iteration are closely related.

Theorem: Let $\{U_n\}_{n=0}^{\infty}$ be the values generated from value iteration, $\{V_n\}_{n=0}^{\infty}$ those from policy iteration. Then if $U_0 = V_0$, $U_n \leq V_n \leq V^*$, for all n .

Proof: We first show that the $\Phi(V) = \max_{\pi} (R_{\pi} + \gamma P_{\pi} V)$ function is monotonic. Indeed, let U, V be s.t. $U \leq V$ and let π be the policy such that

$$\Phi(U) = R_{\pi} + \gamma P_{\pi} U. \text{ Then:}$$

$$\Phi(U) \stackrel{!}{=} R_{\pi} + \gamma P_{\pi} U$$

$$\leq R_{\pi} + \gamma P_{\pi} V$$

$$\leq \max_{\pi} \{R_{\pi} + \gamma P_{\pi} V\}$$

$$= \Phi(V).$$

The main result follows by induction. The base case holds by assumption. Then, suppose $U_n \leq V_n$. By monotonicity of Φ ,

$$U_{n+1} = \Phi(U_n)$$

$$\leq \Phi(V_n)$$

$$\stackrel{!}{=} \max_{\pi} \{R_{\pi} + \gamma P_{\pi} V_n\}.$$

Letting $\tilde{\pi}_{n+1}$ be the policy with $\tilde{\pi}_{n+1} = \arg \max_{\tilde{\pi}} \{R_{\tilde{\pi}} + \gamma P_{\tilde{\pi}} V_n\}$, we ⑦
have

$$\begin{aligned} \Phi(V_n) &= R_{\tilde{\pi}_{n+1}} + \gamma P_{\tilde{\pi}_{n+1}} V_n \\ &\leq R_{\tilde{\pi}_{n+1}} + \gamma P_{\tilde{\pi}_{n+1}} V_{n+1} \end{aligned}$$

which gives the result. $= V_{n+1}$,