MDP  PE

$|S| < \infty$  $|A| < \infty$

# RL Theory[1]: Lecture 3 (Chapter 1)

$s_0 \, a_0 \, r_1 \, S_1 \cdots$
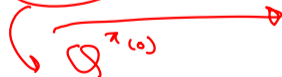
Shivam Garg, RLAI@UAlberta

$V^\pi(s)$

find $V$ for $\pi$

23rd September 2020

PI

$\pi^{(0)}$

$Q^{\pi(0)}$

$\pi^{(1)}(s) = \max_{a'} Q^{\pi(0)}(s, a')$

$Q^{\pi(1)} \cdots$

---

[1]based on https://rltheorybook.github.io/

# Q–Value Iteration Convergence: But what the @#$% does it mean?

- $V^{\pi(k)} \geq V^* - \epsilon\mathbb{1}$ for $k \geq \frac{1}{1-\gamma} \log\left(\frac{2}{\epsilon(1-\gamma)}\right).$ $\Rightarrow$ $k \geq O\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$\log(ab)$
$= \log a + \log b$

$k \geq \left(\frac{1}{1-\gamma}\right)\left[\log\left(\frac{2}{1-\gamma}\right) + \log\left(\frac{1}{\epsilon}\right)\right]$

| Precision | $\epsilon$ | #(iter) = k |
|-----------|------------|-------------|
| 1 | 0.1 | $1.5 \times 10^6$ |
| 2 | 0.01 | $1.7 \times 10^6$ |
| 3 | 0.001 | $1.9 \times 10^6$ |
| 4 | 0.0001 | $2.1 \times 10^6$ |

$\leq O\left(\log\left(\frac{1}{\epsilon}\right)\right)$

- This rate is pretty good! Look at this:

$1 = 0.1$

$\log\left(\frac{1}{\epsilon}\right) = \log(10) = 1$

- This rate is effectively **linear** in the precision of the value function!
- Think of precision as being $\mathbf{log}\ (1/\epsilon)$.

$\log\left(\frac{1}{\epsilon}\right)$

- But this is NOT the sample based update. In practice, we know neither of $r$ and $P$. The rate for the sample based update would be much worse.

# Playing with the Rooftop MDP

- Example calculation from last time:

$$Q^\pi = (1-\gamma)(I - \gamma P^\pi)^{-1} r$$

$$= (1-0.9)\left(\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} - 0.9 \begin{bmatrix} 1 & 0 & 0 \\ 0.5 & 0 & 0.5 \\ 1 & 0 & 0 \end{bmatrix}\right)^{-1} \begin{bmatrix} 0.1 \\ 0.5 \\ 0.2 \end{bmatrix}$$

$$= \begin{bmatrix} Q^\pi(s_1, a_1) \\ Q^\pi(s_1, a_2) \\ Q^\pi(s_2, a_1) \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.1445 \\ 0.11 \end{bmatrix}$$

- Define $\pi_Q(s) := \arg\max_{a \in \mathcal{A}} Q(s, a)$.
- Define $V_Q(s) := \max_{a \in \mathcal{A}} Q(s, a)$.

Figure: (**Rooftop**) MDP.

$$Q^\pi = \begin{bmatrix} Q^\pi(s_1, a_1) \\ Q^\pi(s_1, a_2) \\ Q^\pi(s_2, a_1) \end{bmatrix}$$

$S_1$

$a_1, a_2$

$S_2 - a_1$

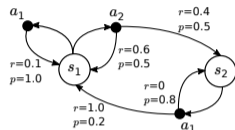# Playing with the Rooftop MDP: Value Iteration



Figure: $V^\pi$



Figure: $Q$–value iteration.

$$Q^* = \begin{bmatrix} 0.29 & 0.32 & 0.27 \end{bmatrix}^\top.$$

# Note about Bellman Operator $\mathcal{T}^\pi$, Bellman Optimality Operator $\mathcal{T}$, and 3 types of $P$ matrices

symbolic

$$Q^\pi = r + \gamma P^\pi Q^\pi \to \text{linear}$$

$Q^\pi$

$$\begin{bmatrix} Q^\pi(s_1 a_1) \end{bmatrix} = \begin{bmatrix} r(s_1 a_1) \end{bmatrix} + \gamma \begin{bmatrix} P(s_1 a_1 | s_1 a_1) \; P(s_1 a_2 | s_1 a_1) \; P(s_2 a_1 | s_1 a_1) \\ \\ \pi \to P^\pi \end{bmatrix} \begin{bmatrix} Q^\pi(s_1 a_1) \\ Q^\pi(s_1 a_2) \\ Q^\pi(s_2 a_1) \end{bmatrix}$$

$s$
$a$
$r(s, a)$
$s' \sim P$
$a' \sim \pi$

matrix

$$Q^\pi(s, a) = r(s, a) + \gamma \sum_{s'} \sum_{a'} P_\pi(s', a' | s, a) Q^\pi(s', a')$$

$?^\pi$     $r$

$$Q^\pi - \gamma P^\pi Q^\pi = r$$

$$\Rightarrow (I - \gamma P^\pi) Q^\pi = r$$

$$\Rightarrow Q^\pi = (I - \gamma P^\pi)^{-1} r$$

# Policy Iteration

- Begin with a policy $\pi^{(0)}$
- Policy Evaluation: Find $Q^{\pi^{(k)}}$; say, by using $Q^\pi = (\phantom{xxx})(I - \gamma P^{\pi^{(k)}})^{-1} r$.
- Policy Improvement: Calculate $\pi^{(k+1)} = \pi_{Q^{\pi^{(k)}}}$.

$\pi^{(0)} \to Q^{\pi^{(0)}} \to \pi^{(1)} \to Q^{\pi^{(1)}} \to \cdots \cdots \to \pi^{*} \qquad \to$ precondition

$? \quad Q^{\pi^{(k)}} \to$

$\pi^{(k+1)}(s) = \arg\max_a Q^{\pi^{(k)}}(s, a)$

- **Convergence:** For $k \geq \frac{1}{1-\gamma} \log\left(1/\epsilon\right)$

$$Q^{\pi_k} \geq Q^* - \epsilon. \qquad \Rightarrow \qquad Q^* - Q^{\pi_k} \leq \epsilon$$

$$k \geq \quad \frac{1}{1-\gamma} \quad \log\left(\frac{1}{\epsilon(1-\gamma)}\right)$$

- I think we need to combine this with the (Singh & Yee, 1994) equation.

# (Part 1) Policy Iteration: Convergence Proof

- $Q^{\pi_{k+1}} \geq \mathcal{T}Q^{\pi_k} \geq Q^{\pi_k}$

$$\mathcal{T}Q = r + \gamma P V_Q$$

$$\begin{bmatrix} V_Q(s_1) \\ V_Q(s_2) \\ \vdots \end{bmatrix}$$

$$V_Q(s) = \max_a Q(s,a)$$

$$\mathcal{T}Q^{\pi_k}(s,a) = r(s,a) + \sum_{s'} p(s'|s,a) \left[ \max_{a'} Q^{\pi_k}(s',a') \right]$$

$$\geq r(s,a) + \sum_{s'} p(s'|s,a) \left[ \sum_{a'} \pi_k(a'|s') Q^{\pi_k}(s',a') \right]$$

$$= Q^{\pi_k}(s,a)$$

# (Part 2) Policy Iteration: Convergence Proof

$$\boxed{Q^{\pi_{k+1}} \geq Q^{\pi_k}}$$

▶ $Q^{\pi_{k+1}} \geq \mathcal{T} Q^{\pi_k} \geq Q^{\pi_k}.$

$$\mathcal{T} Q(s,a) = r(s,a) + \gamma \sum_{s'} p(s'|sa) \max_{a'} Q(s',a')$$

$$Q^{\pi_{k+1}}(s,a) = r(s,a) + \gamma \sum_{s'} p(s'|s,a) \sum_{a'} \pi_{k+1}(a'|s') Q^{\pi_{k+1}}(s',a')$$

$$\geq r(s,a) + \gamma \sum_{s'} p(s'|s,a) \sum_{a'} \pi_{k+1}(a'|s') Q^{\pi_k}(s',a')$$

$$= r(s,a) + \gamma \sum_{s'} p(s'|s,a) \max_{a'} Q^{\pi_k}(s',a')$$

$$= \mathcal{T} Q^{\pi_k}(s,a)$$

$$\pi_{k+1}(s) = \arg\max_a Q^{\pi_k}(s,a)$$

$$\max_{a'} Q^{\pi_k}(s',a')$$

# (Part 3) Policy Iteration: Convergence Proof

- $\|Q^{\pi_{k+1}} - Q^*\|_\infty \leq \gamma \|Q^{\pi_k} - Q^*\|_\infty.$

$$-Q^{\pi_{k+1}} \leq -TQ^{\pi_k} \quad \Big\| \quad Q^* = TQ^*$$

$$\|Q^* - Q^{\pi_{k+1}}\|_\infty \leq \|Q^* - TQ^{\pi_k}\|_\infty$$
$$= \|TQ^* - TQ^{\pi_k}\|_\infty$$
$$\leq \gamma \|Q^* - Q^{\pi_k}\|_\infty \quad (\text{contraction of } T)$$
$$\leq \gamma \left( \gamma \|Q^* - Q^{\pi_{k-1}}\|_\infty \right)$$
$$= \gamma^2 \|Q^* - Q^{\pi_{k-1}}\|_\infty$$
$$\vdots$$
$$\leq \gamma^{k+1} \|Q^* - Q^{\pi_{(0)}}\|_\infty$$

$$\|Q^* - Q^{\pi_{k+1}}\|_\infty$$
$$\leq \gamma^{k+1} \underbrace{\|Q^* - Q^{\pi_{(0)}}\|_\infty}_{\text{fixed}}$$

# (Final Part) Policy Iteration: Convergence Proof

- $$\|Q^* - Q^{\pi_{(k)}}\|_\infty \leq \gamma^k \|Q^* - Q^{\pi_{(0)}}\|_\infty$$

for any $k$

$$= \left(1 - (1-\gamma)\right)^k \|Q^* - Q^{\pi_{(0)}}\|_\infty$$

$$\leq e^{-(1-\gamma)k} \|Q^* - Q^{\pi_{(0)}}\|_\infty \leq \epsilon$$

$$k \geq \frac{1}{1-\gamma} \log\left(\frac{\|Q^* - Q^{\pi_{(0)}}\|_\infty}{\epsilon}\right)$$

$\rightarrow$ atmost this many

$$\|Q^* - Q^{\pi_{(k)}}\|_\infty \leq \epsilon$$

$k \quad O\left(\log(1/\epsilon)\right)$

- How would we show that Generalized Policy Iteration converges?

# Complete Space of Policies

- We define the policy as $\pi : \mathcal{S} \to \Delta(\mathcal{A})$. This is a stationary markov policy.

- A deterministic stationary markov policy would be defined as $\pi : \mathcal{S} \to \mathcal{A}$.

- A general policy (possible non–deterministic, non–stationary, and non–markov) would be defined as $\pi : \mathcal{H} \to \Delta(\mathcal{A})$, where $\mathcal{H}$ is the set of all trajectories.

- Is is even possible to assume the optimal policy to be stationary and deterministic?

# (Part 1) Existence of a Stationary and Deterministic Optimal Policy

- Let $\Pi$ be the set of all non–stationary and randomized policies. There exists a stationary and deterministic policy $\pi$ such that for all states $s \in \mathcal{S}$,

$$V^\pi(s) = \max_{\pi' \in \Pi} V^{\pi'}(s).$$

- Note: I did not understand what randomized means.

# (Part 2) Existence of a Stationary and Deterministic Optimal Policy

-

# (Part 1) Fixed Point of $Q = \mathcal{T}Q$

- Define optimal $Q^*$ by $Q^*(s, a) = \max_{\pi \in \Pi} Q^\pi(s, a)$. Then $Q = Q^*$ if and only if it satisfies

$$Q = \mathcal{T}Q.$$

# (Part 2) Fixed Point of $Q = \mathcal{T}Q$

-

# Summary

- Value iteration

- Policy iteration

- Bellman Operator and Bellman Optimality Operator

- 3 types of transition matrices

- Existence of a stationary and deterministic optimal policy

- Fixed point of $Q = \mathcal{T}Q$

- Did NOT cover the LP formulation (next time)