# RL Theory[1]: Meeting 4 (Chapter 1/2)

Shivam Garg, RLAI@UAlberta

6th October 2020
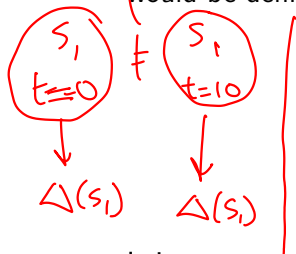
[1]based on https://rltheorybook.github.io/

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

$$V_\phi(s) = \max_a Q(s,a)$$

$$\mathcal{T}Q = r + \gamma P \, V_\phi$$

$$P \quad p(s'|s,s)$$

$$\begin{bmatrix} Q(s,a) \end{bmatrix} \quad \begin{bmatrix} r(s,a) \end{bmatrix} \quad |S||A| \times 1$$

$$\begin{bmatrix} p(s',a'|s,a) & \cdots & p(s',a_s|s,a) \end{bmatrix} \leftarrow P^\pi$$

$$|S \times A| \times |S \times A|$$

$$\mathcal{T}^\pi V = r + \gamma P_\pi^{pred} V$$

$$p(s'|s) \begin{bmatrix} \quad \end{bmatrix} \quad |S| \times |S|$$

# Complete Space of Policies

- We define the policy as $\pi : \mathcal{S} \mapsto \Delta(\mathcal{A})$. 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.

$S_1$
$t=0$

$\neq$

$S_1$
$t=10$

$\triangle(S_1)$    $\triangle(S_1)$

- 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 $s \in \mathcal{S}$,

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

Skip ( didn't understand the proof yet)

▶ 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

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

$$\max_\pi \left( \overset{\pi}{V}(s_1) + \overset{\pi}{V}(s_2) + \overset{\pi}{V}(s_3) \right)$$

$$= \max_\pi V^*(s_1) + \max_\pi V^\pi(s_2) + \cdots$$

$$V^{\pi^*}(s) \geq V^\pi(s) \; \forall \; s \in \mathcal{S}$$

$$\exists \, \pi^*$$

$$Q^*(s,a) := \max_\pi Q^\pi(s,a)$$

$$= \max_\pi \left[ r(s,a) + \gamma \sum_{s'} p(s'|s,a) \overset{\pi}{V}(s') \right]$$

$$= r(s,a) + \gamma \max_\pi \sum_{s'} p(s'|s,a) \overset{\pi}{V}(s')$$

$$= r(s,a) + \gamma \sum_{s'} p(s'|s,a) \underbrace{V^*(s')}_{\downarrow}$$

$$Q^* = \mathcal{T}Q^*$$

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

$$\underset{\text{Optimality}}{\underset{\text{eqn.}}{}} \quad \max_a Q^*(s',a)$$

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

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

$(s,a)$

$$\sum_{s'} p(s'|s,a) \max_{a'} Q(s',a')$$

$P^{\pi_Q} Q$

$$Q = r + \gamma P^{\pi_Q} Q$$

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

for any $\pi'$

$$Q^{\pi'} - Q \leq 0$$

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

$$= (I - \gamma P^{\pi'})^{-1} \left[ I - (I - \gamma P^{\pi'})(I - \gamma P^{\pi_Q})^{-1} \right] r$$

$$= (I - \gamma P^{\pi'})^{-1} \left[ I - \gamma P^{\pi_Q} - (I - \gamma P^{\pi'}) \right] (I - \gamma P^{\pi_Q})^{-1} r$$

$$\underbrace{(I - \gamma P^{\pi'})^{-1} \gamma}_{\geq 0} \underbrace{\left[ P^{\pi'} - P^{\pi_Q} \right] Q}_{\leq 0} \quad \overset{=}{\underset{Q}{}}$$

# Summary

- Value iteration

- Policy iteration

$$\frac{1}{1-x} = 1 + x + x^2 + \cdots$$

$$\left(I - \gamma p^{\pi}\right)^{-1} = I + \gamma p^{\pi} + \gamma^2 (p^{\pi})^2 + \gamma^3 (\gamma^{\pi})^3 + \cdots \geq 0$$

- Bellman Operator and Bellman Optimality Operator

- 3 types of transition matrices

  state-action discounted visitation dist. matrix

  $$p^{\pi} : (s, a) \to (s', a')$$

  $$(p^{\pi})^2 : (s, a) \to (s', a')$$

- Existence of a stationary and deterministic optimal policy

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

  $$\left(p^{\pi} - p^{\pi}_Q\right) Q \leq 0$$

  $$a' \sim \pi_Q$$

  $$\max_a Q(s', a)$$

- Did NOT cover the LP formulation (next time)

$$p^{\pi}_Q Q$$

$$= \sum_{s'} p(s'|s, a) \sum_{a'} \pi_Q(a'|s') Q(s', a')$$

# Sample Complexity

- What is it? Why do we need it?
- Previously, we discussed DP algorithms like value iteration: $Q^{(k+1)} = \mathcal{T}Q^{(k)}$ for $k = 0, 1, 2, \ldots$. Let $\pi^{(k)} = \pi_{Q^{(k)}}$. Then for $k \geq \frac{1}{1-\gamma} \log\left(\frac{2}{\epsilon(1-\gamma)}\right)$,

$$V^{\pi^{(k)}} \geq V^* - \epsilon\mathbb{1}.$$

- This assumes access to the true transition dynamics $P$, which is not available.
- So we address this question: How do these methods perform when we don't have the true $P$?

## Sample Complexity (contd)

▶ We will begin by assuming a naïve model of the environment $\hat{P}$, defined as

$$\hat{P}(s'|s,a) = \frac{\text{count}(s',a,a)}{N}$$

▶ Define $\hat{M}, \hat{V}^{\pi}, \hat{Q}^{\pi}, \hat{Q}^{*}, \hat{\pi}^{*}$.

▶ Then we will see that given the inaccuracy in this model, how accurate our estimates of, say, $\hat{Q}^{\pi}$ can be?

# Sample Complexity for a Naïve Model

- There exists a constant $c$. Let $\epsilon \in \left(0, \frac{1}{1-\gamma}\right)$. If,

$$\# \text{ of samples} \geq \frac{\gamma}{(1-\gamma)^4} \frac{|\mathcal{S}|^2 |\mathcal{A}| \log\left(\frac{c|\mathcal{S}||\mathcal{A}|}{\delta}\right)}{\epsilon^2}$$

then following hold with a probability of greater than $1 - \delta$:

- (Model Accuracy) $\qquad\qquad\qquad \max_{s,a} \|P(\cdot|s,a) - \hat{P}(\cdot|s,a)\|_1 \leq (1-\gamma)^2 \frac{\epsilon}{2}$

- (Uniform Value Accuracy) $\qquad\qquad |Q^\pi - \hat{Q}^\pi\|_\infty \leq \frac{\epsilon}{2} \quad \forall \pi$

- (Near Optimal Planning) $\qquad\qquad \|Q^* - \hat{Q}^*\|_\infty \leq \epsilon$

# Simulation Lemma

- 
$$Q^\pi - \hat{Q}^\pi = \gamma \big( I - \gamma \hat{P}^\pi \big)^{-1} \big( P - \hat{P} \big) V^\pi.$$

# Another Useful Result

- For $x \in |\mathcal{S} \times \mathcal{A}|$

$$\|(I - \gamma \hat{P}^\pi)^{-1} x\|_\infty \leq \frac{\|x\|_\infty}{1 - \gamma}$$