

# RL Theory<sup>1</sup>: Meeting 5 (Chapter 2)

Shivam Garg, RLAI@UAlberta

20th October 2020

---

<sup>1</sup>based on <https://rltheorybook.github.io/>

# Sample Complexity

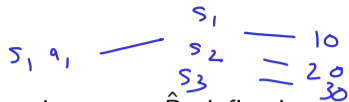
*Q-value iteration*

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

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

- ▶ 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)

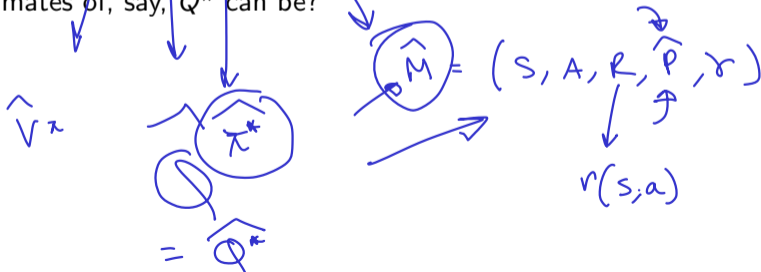


- ▶ We will begin by assuming a naïve model of the environment  $\hat{P}$ , defined as

$$s' \sim P(\cdot | s, a)$$

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

- ▶ 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  $\epsilon \in (0, \frac{1}{1-\gamma})$ . If,

$$N(s, a) = N$$

# (s, a) ↗

$|S||A|N$

$\mathcal{T} \xrightarrow{\hat{P}} \mathcal{Q}$

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

$s_1, a_1 - 100$   
 $s_1, a_2 - 10$

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

- ▶ (Model Accuracy)
- ▶ (Uniform Value Accuracy)
- ▶ (Near Optimal Planning)

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


$$\|Q^\pi - \hat{Q}^\pi\|_\infty \leq \frac{\epsilon}{2} \quad \forall \pi$$

$$\|Q^* - \hat{Q}^*\|_\infty \leq \epsilon$$

$\leq \epsilon'/2$   
 $\epsilon' = (1-\gamma)\epsilon$

# Hoeffding's Inequality

▶


$$\mathbb{P} \left( \left| \mathbb{E}[X] - \frac{\sum_{i=1}^N X_i}{N} \right| \leq (b_+ - b_-) \sqrt{\frac{\ln(2/\delta)}{2N}} \right) \geq 1 - \delta.$$

*scalar*

$X_1, X_2, X_3, \dots \sim P$

$X_i \in [b_-, b_+]$

↓                  ↓  
1                    6

Hoeffding

# Union Bound

$$A_1, A_2, A_3, \dots \in \mathcal{A}$$

$(s, a) \rightarrow P(\cup_i A_i) \leq \sum P(A_i)$

$$\mathbb{P}(|A_i| \leq c(\delta)) \geq 1 - \delta \Rightarrow \mathbb{P}\left(\max_i |A_i| \leq c(\delta)\right) \geq 1 - \underbrace{|\mathcal{A}|\delta}_{\delta'}$$

$\hat{P}(s'/s, a) \leq c$

$$1 - P(|A_i| \leq c) \leq \delta$$

$$P(|A_i| \geq c) \leq \delta$$

$$1 - P(\cup_i (|A_i| \geq c)) \geq 1 - |\mathcal{A}|\delta$$

$$P\left(\underbrace{\bigcup_i |A_i| \geq c}_{\text{red underline}}\right) \geq 1 - |\mathcal{A}|\delta$$

$$\bigcap |A_i| \leq c \Rightarrow \max_i |A_i| \leq c$$

# Proof of Claim 1: (Step 1)

$\delta \ll 1$

$$P\left(\|P(\cdot|s, a) - \hat{P}(\cdot|s, a)\|_1 \leq c \sqrt{\frac{|S| \log(1/\delta)}{m}}\right) \geq 1 - \delta$$

$q = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$   
↓  
 $d_{x,1}$

$$P\left(\|\hat{q} - \bar{q}\|_1 \leq \sqrt{|S|} \left(\frac{1}{\sqrt{m}} + \epsilon\right)\right) \geq 1 - e^{-N\epsilon^2}$$

$\frac{\sum q_i}{N}$

$E[q]$

$$\epsilon \leq \sqrt{\frac{1}{m} \log(1/\delta)}$$

$\geq 1 - \delta$

$\epsilon \in (0, \frac{1}{1+\delta})$

$\delta \approx 1$

$$\left(\frac{\sqrt{|S|}}{\sqrt{m}} + \sqrt{\frac{|S|}{m} \log(1/\delta)}\right) \leq C \sqrt{\frac{|S|}{m} \log(1/\delta)}$$

# Proof of Claim 1: (Step 2)

$\# \text{ samples} \geq \frac{\gamma}{(1-\gamma)^4} \frac{|S|^2 |\mathcal{A}| \log\left(\frac{c|S||\mathcal{A}|}{\delta}\right)}{\epsilon^2} \Rightarrow \max_{s,a} \|P(\cdot|s,a) - \hat{P}(\cdot|s,a)\|_1 \leq (1-\gamma)^2 \frac{\epsilon}{2}$

$\mathbb{P}\left(\max_{s,a} \|p(\cdot|s,a) - \hat{p}(\cdot|s,a)\| \leq c \sqrt{\frac{|S| \log\left(\frac{|S||\mathcal{A}|}{\delta}\right)}{m}}\right) \geq 1-\delta$

$\frac{1}{m} = \frac{(1-\gamma)^4 \epsilon^2}{4 |S| \log\left(\frac{|S||\mathcal{A}|}{\delta}\right)} \Rightarrow |S||\mathcal{A}| m \geq \frac{4 |\mathcal{A}| |S|^2 \log\left(\frac{|S||\mathcal{A}|}{\delta}\right)}{(1-\gamma)^4 \epsilon^2 c^2}$



## Proof of Claim 1: (Step 3)



$$\# \text{ samples} \geq \frac{\gamma}{(1-\gamma)^4} \frac{|\mathcal{S}|^2 |\mathcal{A}| \log\left(\frac{c|\mathcal{S}||\mathcal{A}|}{\delta}\right)}{\epsilon^2} \Rightarrow \max_{s,a} \|P(\cdot|s,a) - \hat{P}(\cdot|s,a)\|_1 \leq (1-\gamma)^2 \frac{\epsilon}{2}$$



# Simulation Lemma

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

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

$$\underbrace{(I - \gamma P^\pi)^{-1}} r - \underbrace{(I - \gamma \hat{P}^\pi)^{-1}} r$$

$$\gamma (I - \gamma \hat{P}^\pi)^{-1} \left[ P^\pi - \hat{P}^\pi \right] \underbrace{(I - \gamma P^\pi)^{-1} r}_{Q^\pi}$$

$$= \sum_{a'} \pi(a' | s') \sum_{s'} P(s' | s, a) \underbrace{[P^\pi - \hat{P}^\pi] Q^\pi}_{(P - \hat{P}) V} = \sum_{s'} P(s' | s, a) \sum_{a'} \pi(a' | s') \underbrace{Q^\pi(s', a')}_{V^\pi(s')}$$

## Another Useful Result

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

$$\frac{1}{1-\gamma} = 1 + \gamma + \gamma^2 + \dots$$

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

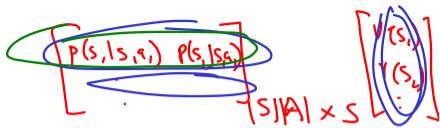
$$\begin{aligned} & \| (I + \gamma \hat{P}^\pi + \gamma^2 \hat{P}^\pi + \dots) x \|_\infty \\ & \leq \|x\|_\infty + \gamma \| \hat{P}^\pi x \|_\infty + \dots \\ & \leq (1 + \gamma + \gamma^2 + \dots) \|x\|_\infty \end{aligned}$$

Proof of Claim 2

Hölder's

$r \in [0, 1]$

$\left[ \begin{array}{c} \vdots \\ p \\ \vdots \end{array} \right] \rightarrow \sum_i |p_i| |q_i| \leq \|P\|_1 \|q\|_\infty \quad \|\hat{Q}^\pi - Q^\pi\|_\infty \leq \epsilon/2.$



$$\|\hat{Q}^\pi - Q^\pi\|_\infty = \gamma \left\| \underbrace{(\mathbb{I} - \gamma \hat{P}^\pi)^{-1}}_{\text{matrix}} \underbrace{(P - \hat{P}) v^\pi}_{\text{vector}} \right\|_\infty$$

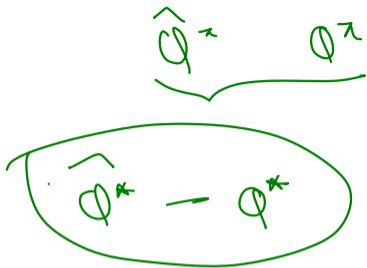
$$\leq \frac{\gamma}{1-\gamma} \underbrace{\| (P - \hat{P}) v^\pi \|_\infty}_{\text{vector}}$$

$P(\cdot | s, a) \quad v^\pi$

$$\begin{aligned} & \max_{s, a} \left| \sum_{s'} (p(s' | s, a) - \hat{p}(s' | s, a)) v^\pi(s') \right| \\ & \leq \max_{s, a} \sum_{s'} |p(s' | s, a) - \hat{p}(s' | s, a)| |v^\pi(s')| \\ & \max_{s, a} \underbrace{\| p(\cdot | s, a) - \hat{p}(\cdot | s, a) \|_1}_{\text{matrix}} \|v^\pi\|_\infty \rightarrow \frac{1}{1-\gamma} \end{aligned}$$

$\frac{\gamma}{(1-\gamma)^2} (1-\gamma)^2 \epsilon^2/2$

$\frac{\gamma}{1-\gamma} (\epsilon/2) \leq \frac{\gamma}{1-\gamma}$



$$\frac{\sum X_i}{n} \left| \hat{P}(s'|s,a) \right| X_i = \mathbb{I}(s_i = s')$$

Diagram illustrating the relationship between the sample mean  $\frac{\sum X_i}{n}$  and the indicator function  $\mathbb{I}(s_i = s')$ . A red arrow points from  $\frac{\sum X_i}{n}$  to  $\mathbb{I}(s_i = s')$ . Another red arrow points from  $\mathbb{I}(s_i = s')$  to  $E[X_i]$ .

$$\| (P - \hat{P}) v^* \|_\infty = \max_{s,a} \left| \sum_{s'} P(s'|s,a) v^*(s') - \sum_{s'} \hat{P}(s'|s,a) v^*(s') \right|$$

$$= \frac{1}{1-\gamma} \sqrt{\frac{2}{n} \log \left( \frac{2|S||A|}{\delta} \right)}$$