

CS292F StatRL Lecture 3

MDP with a generative model

Instructor: Yu-Xiang Wang

Spring 2021

UC Santa Barbara

Recap: Markov Decision processes (MDP)

- Infinite horizon / discounted setting

$$\mathcal{M}(\mathcal{S}, \mathcal{A}, P, r, \gamma, \mu)$$

Transition kernel: $P: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ i.e. $P(s'|s, a)$

(Expected) reward function: $r: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R} / [0, R_{\max}]$ $\mathbb{E}[R_t | S_t=s, A_t=a] =: r(s, a)$

Initial state distribution $\mu \in \Delta(\mathcal{S})$

Discounting factor: γ

Recap: Reward function and Value functions

$(S_1, A_1, R_1, \dots, S_T, A_T, R_T, \dots)$

- Immediate reward function $r(s,a)$

- **expected immediate** reward

$$r(s, a) = \mathbb{E}[R_1 | S_1 = s, A_1 = a]$$

$$r^\pi(s) = \mathbb{E}_{a \sim \pi(a|s)}[R_1 | S_1 = s]$$

$S = |S|$
 $A = |A|$

- state value function: $V^\pi(s)$

- **expected long-term** return when starting in s and following π

$$V^\pi(s) = \mathbb{E}_\pi[R_1 + \gamma R_2 + \dots + \gamma^{t-1} R_t + \dots | S_1 = s]$$

- state-action value function: $Q^\pi(s,a)$

- **expected long-term** return when starting in s , performing a , and following π

$$Q^\pi(s, a) = \mathbb{E}_\pi[R_1 + \gamma R_2 + \dots + \gamma^{t-1} R_t + \dots | S_1 = s, A_1 = a]$$

Recap: Bellman equations

$$Q = r + \gamma P Q$$

$$V^\pi = r^\pi + \gamma P^\pi V^\pi$$

$$Q^\pi = r + \gamma P V^\pi$$

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

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

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



discounted-occupancy measures

$$V^\pi = \mu + \gamma (P^\pi)^T V^\pi$$

$$\tilde{V}^\pi = \mu^\pi + \gamma (P^\pi)^T \tilde{V}^\pi$$

$$\tilde{V}_{(s_a)}^\pi = V^\pi(s) \cdot \pi(a|s)$$

$$V^\pi = (I - \gamma (P^\pi)^T)^{-1} \mu$$

$$\tilde{V}^\pi = (I - \gamma (\tilde{P}^\pi)^T)^{-1} \mu^\pi$$

$$\mu_{(s_a)}^\pi = \mu^\pi(s) \pi(a|s)$$

Recap: Duality and LP-formulation

- Primal LP:

$$\min \sum_s \mu(s)V(s)$$

$$\text{subject to } V(s) \geq r(s, a) + \gamma \sum_{s'} P(s'|s, a)V(s') \quad \forall a \in \mathcal{A}, s \in \mathcal{S}$$

V^* is the optimal solution of the primal LP
 fix π , V^π , the obj func = $\langle V^\pi, \mu \rangle = V^\pi(\mu)$
 the value of obj π

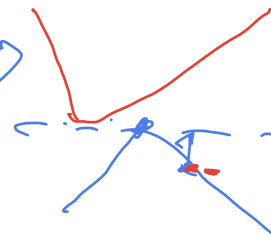
- Dual LP:

$$\max_\nu \sum_{s,a} \nu(s, a)r(s, a)$$

$$\text{subject to } \nu \geq 0$$

$$\sum_z \nu(s, a) = \mu(s) + \gamma \sum_{s', a'} P(s|s', a')\nu(s', a')$$

ν^* is the optimal solution to the dual LP
 for a fixed π , $V^\pi(s, a)$ is feasible
 for dual LP



Recall: Bellman optimality equation and a stationary and deterministic optimal policy

$$V^*(s) = \max_a \sum_{s'} P(s'|s, a) [r(s, a, s') + \gamma V^*(s')]$$

$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

Handwritten notes: $Q^*(s, a) = r(s, a) + \sum_{s'} P(s'|s, a) V^*(s')$, $V^*(s) = \max_a Q^*(s, a)$, $Q^*(s, a) = r(s, a) + \max_{a'} \sum_{s'} P(s'|s, a') Q^*(s', a')$

Value iterations (VI) aim at finding the fixed point by recursively applying the Bellman (optimality) operator.

Lemma 1. The Bellman operator is a γ -contraction.

For any two vectors $Q, Q' \in \mathbb{R}^{|S| \times |A|}$,

$$\|TQ - TQ'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

Recap: computational complexity

$$\log \frac{1}{\epsilon}$$

Repeat $\left\{ \begin{array}{l} \textcircled{1} \text{ Policy Evaluation of } \pi_t : \text{Solve Bellman equation} \\ \textcircled{2} \text{ Policy improvement } \pi_{t+1} = \arg \max_{\pi} Q^{\pi_t} \end{array} \right.$

	Value Iteration	Policy Iteration	LP-Algorithms
Poly?	$ \mathcal{S} ^2 \mathcal{A} \frac{L(P, r, \gamma) \log \frac{1}{1-\gamma}}{1-\gamma}$	$(\mathcal{S} ^3 + \mathcal{S} ^2 \mathcal{A}) \frac{L(P, r, \gamma) \log \frac{1}{1-\gamma}}{1-\gamma}$	$ \mathcal{S} ^3 \mathcal{A} L(P, r, \gamma)$
Strongly Poly?	X	$(\mathcal{S} ^3 + \mathcal{S} ^2 \mathcal{A}) \cdot \min \left\{ \frac{ \mathcal{A} ^{ \mathcal{S} }}{ \mathcal{S} }, \frac{ \mathcal{S} ^2 \mathcal{A} \log \frac{ \mathcal{S} ^2}{1-\gamma}}{1-\gamma} \right\}$	$ \mathcal{S} ^4 \mathcal{A} ^4 \log \frac{ \mathcal{S} }{1-\gamma}$

A trivial lower bound: $\Omega(SA)$ needed to store the Q^* function.

$$\underline{S \log |\mathcal{A}|}$$

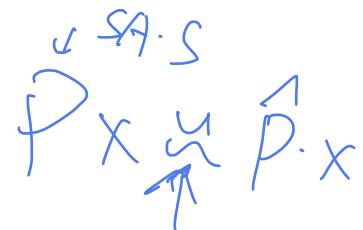
Question: If we allow **randomness**, can we further improve the computational complexity?

- Large MDPs

- Backgammon: 10^{20}
- Chess: 10^{47}
- Game of Go: 10^{174}

- The transition kernel requires S^2A parameters to describe, and to apply.

- VI, PI, LP all depends at least S^2A



- **What if we can sample transition in $O(1)$?**

Access to a simulator or a generative model $\underline{S'} \sim P(\cdot | s, a)$

- Popularized by [Kakade \(2003\)](#)
- Examples when this is a meaningful model:
 - Games
 - Robotics simulation
 - RL for Science
- Not the most realistic if
 - The simulator is a crude approximation of the world
 - You cannot take a snapshot and restart

How many generative model oracle calls do we need to obtain an ϵ -optimal policy?

- (oracle) computational complexity
 - Assume $O(1)$ time to draw sample $S' \sim P(\cdot | s, a)$
- But also can be viewed as a simplified version of the sample complexity of RL
 - without worrying about exploration.
 - Let's get N samples **for each** (s,a) pairs.
 - How is N related to ϵ

How are we using the simulator?

- We will consider the dumbest way of using it
 - Sampling N rounds. Each round go over each (s,a) pair.

$$\text{total} = N \cdot |S| \cdot |A| \quad \overline{S+A}$$

- A total of NSA oracle calls.

- We have N samples for each SA, but often $N \ll S$

$$\frac{S^2 A \log\left(\frac{1}{\epsilon(\delta)^2}\right)}{1-\delta}$$

- It is possible to do better than this, but not in the worst case, so we will study this algorithm first.

Plug-in estimator of P

$$\hat{P}(s'|s, a) = \frac{\text{count}(s', s, a)}{N} \text{ where } \text{count}(s', s, a) = \sum_{i=1}^N \mathbf{1}(S'_{i,s,a} = s').$$

for each s, a
we sample i.i.d.
 $S'_{1,s,a}, \dots, S'_{N,s,a}$
 $\sim \text{Ber}(P(s'|s,a))$

- How many parameters does P have?

$S^2 \cdot A$

- Often in large MDP, $N \ll S$

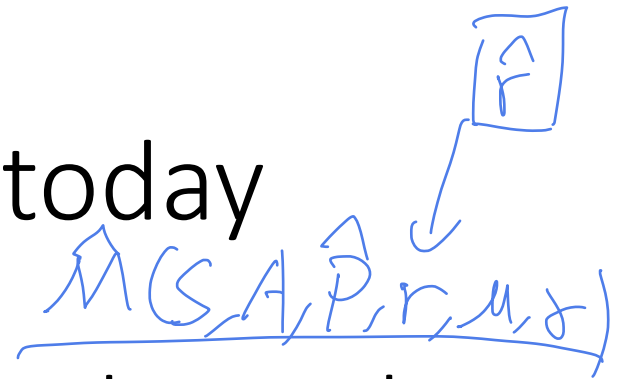
$N \cdot S \cdot A$

$\hat{P}(s'|s,a) = 0$ for many s'

Key question of interest:

Do we need to estimate P accurately to obtain near optimal policies?

Outline of the lecture today



- Simulation Lemma and model-based approach
- Review of statistical tools we need:
 - Hoeffding's inequality
 - Bernstein inequality
 - McDiarmid's inequality
- Sample complexity bounds

Model-based approach

Q^* is optimal \hat{M}

- Approximate MDP $\hat{M} = (S, A, \hat{P}, r, \gamma, \delta)$
 - Run VI, PI on the approximate MDP

$$Q^* \leftarrow VI(\hat{M})$$

$$Q^* = \underset{Q}{\operatorname{argmax}} Q(S, a)$$

Q^* is the optimal Value function of \hat{M}

- From uniform convergence to suboptimality bound

$$\underbrace{V^{\pi^*} - V^{\hat{\pi}^*}}_{\leq \epsilon} = \underbrace{V^{\pi^*} - V^{\pi^*}}_{\leq \epsilon} + \underbrace{V^{\pi^*} - V^{\hat{\pi}^*}}_{\leq 0} + \underbrace{V^{\hat{\pi}^*} - V^{\hat{\pi}^*}}_{\leq \epsilon} \leq 2\epsilon$$

Uniform Convergence

$$\sup_{\pi \in \Pi} \|V^{\pi} - V^{\hat{\pi}}\|_{\infty} \leq \epsilon$$

Computational complexity of the model-based approach

- To construct the approximate transition kernel

for each (S, a) sample N

time $N \cdot SA$

space $\# \text{ of nonzeros } : SA$
 $\leq \min N \cdot SA, S^2 A$

- To compute empirically optimal policy
 - via value iteration

$$\left. \begin{array}{l} \min N \cdot SA \\ S^2 A \end{array} \right\} \cdot \frac{\log \frac{1}{\epsilon} (1+\delta)^2}{1-\delta}$$

Attempt 1: Simulation Lemma (Kearns and Singh, 2002)

Lemma 2.2. (Simulation Lemma) For all π we have that:

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

- Proof using closed-form solution for Q

$$\begin{aligned}
 Q^\pi - \hat{Q}^\pi &= (I - \gamma P^\pi)^{-1} r - (I - \gamma \hat{P}^\pi)^{-1} r \\
 &= (I - \gamma \hat{P}^\pi)^{-1} ((I - \gamma \hat{P}^\pi) - (I - \gamma P^\pi)) Q^\pi \\
 &= \gamma (I - \gamma \hat{P}^\pi)^{-1} (P^\pi - \hat{P}^\pi) Q^\pi \\
 &= \gamma (I - \gamma \hat{P}^\pi)^{-1} (P - \hat{P}) V^\pi
 \end{aligned}$$

$$P^\pi(s', a' | s, a) = P(s' | s, a) \cdot \pi(a' | s')$$

Uniform convergence via the Simulation Lemma

$$\sup_{\pi} \|Q^{\pi} - \hat{Q}^{\pi}\|_{\infty} = \|\gamma(I - \gamma\hat{P}^{\pi})^{-1}(P - \hat{P})V^{\pi}\|_{\infty} \leq \frac{\gamma}{1-\gamma} \|(P - \hat{P})V^{\pi}\|_{\infty} \quad (1)$$

$$\leq \frac{\gamma}{1-\gamma} \left(\max_{s,a} \|P(\cdot|s,a) - \hat{P}(\cdot|s,a)\|_1 \right) \|V^{\pi}\|_{\infty} \quad (2)$$

$$\leq \frac{\gamma}{(1-\gamma)^2} \max_{s,a} \|P(\cdot|s,a) - \hat{P}(\cdot|s,a)\|_1 \leq \frac{\gamma}{1-\gamma} \quad (3)$$

- We proved (1) when we prove the invertability in Lecture 2

$$\begin{aligned} (1): \|(I - \gamma\hat{P}^{\pi})^{-1}x\|_{\infty} &\leq \frac{1}{1-\gamma} \|x\|_{\infty} & \|(I - \gamma\hat{P}^{\pi})y\|_{\infty} &\geq \|y\|_{\infty} - \gamma\|\hat{P}^{\pi}y\|_{\infty} \\ & & &\geq \|y\|_{\infty}(1-\gamma) \\ \text{and} \quad \|(I - \gamma\hat{P}^{\pi})y\|_{\infty} &\geq (1-\gamma)\|y\|_{\infty} \end{aligned}$$

- Key observation: RHS doesn't depend on the policy.**

All (you need to know) about ^{the intuition of} Statistics in one slide, two theorems.

- Statistics is about using **samples** from a distribution to infer the properties of the distribution itself (**population**)

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



- Law of large number
 - Average \rightarrow Mean

$$\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow[n \rightarrow \infty]{\text{a.s.}} \mathbb{E}[X_1]$$

- Central limit theorem
 - The rate is sqrt(1/n)

$$\sqrt{n} \cdot \left(\frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E}[X_1] \right) \xrightarrow[\text{Law}]{\text{in}} N(0, \text{Var}(X_1))$$

Concentration inequalities --- finite-sample bounds of LLN and CLT

$$\begin{aligned} & \text{any } X \leq B \\ & \text{Var}(X) \leq \frac{B^2}{4} \end{aligned}$$

- **Hoeffding's inequality:** Assume X_1, \dots, X_n are independent and their support bounded: $P(a_i \leq X_i \leq b_i) = 1$

$$S_n = X_1 + \dots + X_n$$

$$P\left(\frac{S_n - \mathbb{E}[S_n]}{n} \geq \frac{t}{n}\right) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

- Easy version, if $0 < X_i < B$, with probability $1 - \delta$:



$$\begin{aligned} & e^{-\frac{2nB^2 \log(2/\delta)}{t^2}} \\ & \text{or } B^2 \\ & = e^{-\log(2/\delta)} = \frac{\delta}{2} \end{aligned}$$

$$|\bar{X} - \mathbb{E}[\bar{X}]| \leq \sqrt{\frac{B^2}{2n} \log(2/\delta)}$$

$$\frac{t}{n} = \sqrt{\frac{B^2}{2n} \log(2/\delta)} \quad \text{then } t = \sqrt{\frac{2nB^2 \log(2/\delta)}{2}}$$

Concentration inequalities --- finite-sample bounds of LLN and CLT

$Y_1 \dots Y_n$ $X_i = Y_i - \mathbb{E}Y_i$

- **Bernstein inequality:** Assume X_1, \dots, X_n are independent, **zero-mean**, and their absolute value bounded by M , then

$$\mathbb{P} \left(\sum_{i=1}^n X_i \geq t \right) \leq \exp \left(- \frac{\frac{1}{2}t^2}{\underbrace{\sum_{i=1}^n \mathbb{E}[X_i^2]}_{\text{Var}(X_i)} + \frac{1}{3}Mt} \right).$$

- Easy version for the iid case, with probability 1-δ:

$$|\bar{X} - \mathbb{E}[X_1]| \leq \underbrace{\sqrt{\frac{2\text{Var}[X_1]}{n} \log(2/\delta)}}_{\substack{\uparrow \\ \text{CLT} \\ O(\frac{1}{\sqrt{n}})}} + \underbrace{\frac{2M \log(2/\delta)}{3n}}_{\substack{\text{circled} \\ O(\frac{1}{n})}}$$

A generalization of Hoeffding's inequality to McDiarmid's Inequality

McDiarmid's inequality: Assume X_1, \dots, X_n are independent, and function f satisfies the following

Coordinatewise Uniform Stability condition:

$$\sup_{x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n} |f(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n) - f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)| \leq c_i.$$

Then we have:

$$\left\| \frac{1}{N} \sum x_i - P \right\|_1 - \left\| \frac{1}{N} \sum x_i - x_j + x'_j - P \right\|_1 \leq \|x_j + x'_j\|_1 \leq 2$$

$$P(f(X_1, X_2, \dots, X_n) - \mathbb{E}[f(X_1, X_2, \dots, X_n)] \geq \epsilon) \leq \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^n c_i^2}\right),$$

$f(\hat{P}, P) = \|\hat{P} - P\|_1$
 $X_i = \sum_{(s_1, \dots, s_n)} e_{s_1, \dots, s_n}$
 $\hat{P} = \frac{1}{N} \sum_i e_{s_1, \dots, s_n}$
 $X_i = \vec{e}_{s_1, \dots, s_n} - P(\cdot | s_n)$

High probability bound for estimating probability distribution in L1

- Apply McDiarmid inequality

$$f(x_1, \dots, x_n) - E f(x_1, \dots, x_n)$$

$$\left| \|\hat{P} - P\|_1 - E[\|\hat{P} - P\|_1] \right| < \sqrt{\frac{c^2 \log^2 \frac{2}{\delta}}{2n}} = \sqrt{\frac{2 \log^2 \frac{2}{\delta}}{n}}$$

$n = N$

- Calculate the expectation

$$\|\hat{P} - P\|_1 \leq \left| \|\hat{P} - P\|_1 - E[\|\hat{P} - P\|_1] \right| + E[\|\hat{P} - P\|_1] \leq \sqrt{\frac{2 \log^2 \frac{2}{\delta}}{n}} + \sqrt{\frac{1}{n}}$$

$$E[\|\hat{P} - P\|_1] \leq \sqrt{1} E[\|\hat{P} - P\|_2] \leq \sqrt{1} \sqrt{E[\|\hat{P} - P\|_2^2]} = \sqrt{1} \sqrt{\sum_s E(\hat{P}(s) - P(s))^2}$$

↑ Jensen's inequality

$$= \sqrt{1} \sqrt{\sum_s \frac{P(s)(1-P(s))}{n}}$$

$$\leq \sqrt{\frac{1}{n}}$$

$$\frac{\sqrt{\sum x_i^2}}{\sqrt{\sum x_i^2}} = 1$$

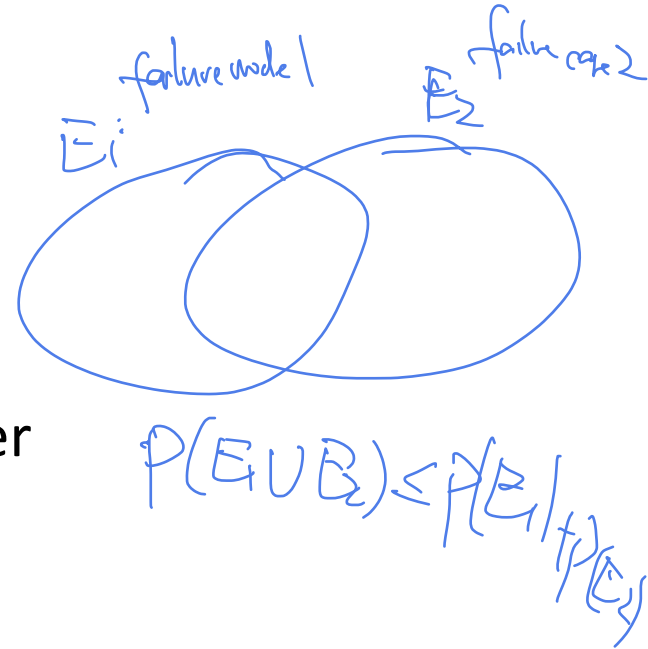
$$\| \cdot \|_1 \leq \sqrt{d} \| \cdot \|_2$$

The “Union bound” trick

- “Union bound”: For a countable sequence of events A_1, A_2, A_3, \dots

$$\mathbb{P}\left(\bigcup_i A_i\right) \leq \sum_i \mathbb{P}(A_i).$$

- Typical use case:
 - Bound low-probability events together



Continue with the uniform convergence via the "Simulation Lemma"

$$\|Q^\pi - \hat{Q}^\pi\|_\infty = \|\gamma(I - \gamma\hat{P}^\pi)^{-1}(P - \hat{P})V^\pi\|_\infty \leq \frac{\gamma}{1-\gamma} \|(P - \hat{P})V^\pi\|_\infty \quad (1)$$

$$\leq \frac{\gamma}{1-\gamma} \left(\max_{s,a} \|P(\cdot|s,a) - \hat{P}(\cdot|s,a)\|_1 \right) \|V^\pi\|_\infty \quad (2)$$

$$\leq \frac{\gamma}{(1-\gamma)^2} \max_{s,a} \|P(\cdot|s,a) - \hat{P}(\cdot|s,a)\|_1 \quad (3)$$

$\leq \frac{\delta}{(1-\gamma)^2} \left(\sqrt{\frac{2 \log \frac{2(S|A)}{\delta}}{\delta}} + \sqrt{S} \right)$
w.p. $1-\delta$ $\|S|A\| \cdot \delta$

$\sup_{s,a} \|\hat{P}(\cdot|s,a) - P\|_1 \leq \sqrt{\frac{2 \log \frac{2}{\delta}}{N}} + \sqrt{\frac{S}{N}}$

$\delta = |S| |A| \delta$

$\leq \epsilon$

if suffices $N \geq \frac{2\gamma^2 (\log \frac{2S|A}{\epsilon} + S)}{(1-\gamma)^2 \epsilon}$

Summary of Attempt 1: “simulation lemma” + uniform convergence

- Sample complexity

$$N \geq \frac{2\sigma^2 \left(\log\left(\frac{SA}{\delta}\right) + S \right)}{(1-\gamma)^4 \epsilon^2}$$

- Computational complexity

$$\mathcal{O}\left(N S A + \underbrace{S^2 A}_{\max(N, S)} \log(1/\delta) \right)$$
$$S^2 A$$

Exercise: Try the alternatives

- Try applying Hoeffding's inequality coordinatewise, then union bound over s'
 - Could you recover the same bound?
- Try applying Bernstein's inequality coordinatewise, then union bound over s'
 - Do you need additional assumptions to get the same bound?

Attempt 2: Bounding the value function instead

- Recall:

Lemma 1.11 AJKS (Q-error amplification):

$$\underline{V^{\pi_Q}} \geq V^* - \frac{2\|Q - Q^*\|_\infty}{1 - \gamma} \mathbb{1}.$$

- If we can bound $\|\hat{Q}^* - Q^*\|_\infty$ with an error independent to S , then we can improve the previous bound

$Q \in \mathbb{R}^{SA}$ $\pi_Q = \underset{a}{\operatorname{Argmax}} Q(s,a)$

Bounding the value function

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

- Key lemma: $\|Q^* - \hat{Q}^*\|_\infty \leq \frac{\gamma}{1-\gamma} \|(P - \hat{P})V^*\|_\infty$

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

- Proof: Use the contraction of Bellman (optimality) operator.

$$\begin{aligned} \|Q^* - \hat{Q}^*\|_\infty &= \|Q^* - \mathcal{T}Q^* + \mathcal{T}Q^* - \mathcal{T}\hat{Q}^*\|_\infty \\ &\stackrel{\text{Def. of } \mathcal{T}}{\leq} \|Q^* - \mathcal{T}Q^*\|_\infty + \|\mathcal{T}Q^* - \mathcal{T}\hat{Q}^*\|_\infty \\ &\leq \|r + \gamma P V^* - (r + \gamma \hat{P} \frac{\max_{a \in \mathcal{A}(s)} Q^*(s,a)}{V^*})\|_\infty + \gamma \|Q^* - \hat{Q}^*\|_\infty \leftarrow \text{Bellman Contraction} \\ &= \|\gamma(P - \hat{P}) \cdot V^*\|_\infty + \gamma \|Q^* - \hat{Q}^*\|_\infty \end{aligned}$$

The key trick for knocking off an S factor is the following:

- Key lemma: $\|Q^* - \hat{Q}^*\|_\infty \leq \frac{\gamma}{1-\gamma} \|(P - \hat{P})V^*\|_\infty$

Handwritten notes:

$$\begin{aligned} P \cdot V^* &= \frac{1}{N} \sum e_{s_i}^T V^* \\ \hat{P} &= \frac{1}{N} \sum e_{s_i} \end{aligned}$$

$$\begin{aligned} \|(P - \hat{P})V^*\|_\infty &= \max_{s,a} \left| E_{s' \sim P(\cdot|s,a)}[V^*(s')] - E_{s' \sim \hat{P}(\cdot|s,a)}[V^*(s')] \right| \\ &= \max_{s,a} \left| E_{s' \sim P(\cdot|s,a)}[V^*(s')] - \frac{1}{N} \sum_{i=1}^N V^*(S'_{i,s,a}) \right| \end{aligned}$$

Handwritten note: $E[V^*]$ with an arrow pointing to the expectation term in the equation above.

Apply Hoeffding's inequality!

u.p. $1-\delta$

Handwritten note: $0 \leq V^* \leq \frac{1}{1-\gamma}$

$$\left| \frac{1}{N} \sum V^*(S'_{i,s,a}) - E_{s'}[V^*(s')] \right| \leq \sqrt{\frac{B^2 \log \frac{2}{\delta}}{2N}} = \frac{1}{1-\gamma} \sqrt{\frac{\log \frac{2}{\delta}}{2N}}$$

Summary of Attempt 2: "Q-amplification" + Bellman operator

- Sample complexity

$$V^* - V^{\frac{1}{\gamma^*}} \leq \frac{2 \|Q^* - Q\|_{\infty}}{1 - \gamma} \leq \frac{2\delta}{(1-\gamma)^3} \sqrt{\frac{\log(2SA)}{\epsilon}} \quad \text{w.p. } \frac{\epsilon}{1-\delta}$$

$2N = \epsilon$

- Computational complexity

$$N \cdot SA + N \cdot SA \cdot \log(\text{\# of iterations for VI})$$

$$N \geq \frac{2\delta \log \frac{2SA}{\epsilon}}{(1-\gamma)^3 \epsilon^2}$$

Optimal sample complexity (Azar et al., 2013)

$$N = \Theta \left(\frac{1}{(1 - \gamma)^3} \frac{\log(cSA/\delta)}{\epsilon^2} \right)$$

Recent literature:

- (Sidford et al, 2018) A variance reduced approx. value iteration-based approach for $\epsilon < 1$
- (Agarwal et al., 2019) Proven the same for model-based approach for $\epsilon < \sqrt{1/(1-\gamma)}$
- (Li et al., 2020) optimal rates for all values of $\epsilon < 1/(1-\gamma)$ for a perturbed model-based approach.
- (Yin, Bai, W., 2020) optimal rates for the finite horizon case with model-based plug-in method. $\epsilon < \sqrt{H}$ $H \approx \frac{1}{1-\gamma}$
- (Yin, Bai, W., 2021) double variance reduction, all values of $\epsilon < H$, finite horizon case (and $\epsilon < 1/(1-\gamma)$ for the infinite horizon case too)

It remains an open problem whether model-based plug-in is optimal for all ϵ

References on the minimax sample complexity

Azar, M. G., Munos, R., & Kappen, H. J. (2013). Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3), 325-349.

Sidford, A., Wang, M., Wu, X., Yang, L. F., & Ye, Y. (2018). Near-optimal time and sample complexities for solving Markov decision processes with a generative model. *Advances in Neural Information Processing Systems*

Agarwal, A., Kakade, S., & Yang, L. F. (2020). Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory* (pp. 67-83). PMLR.

Li, G., Wei, Y., Chi, Y., Gu, Y., & Chen, Y. (2020). Breaking the sample size barrier in model-based reinforcement learning with a generative model. *Advances in Neural Information Processing Systems*, 33.

Yin, M., Bai, Y., & Wang, Y. X. (2020). Near optimal provable uniform convergence in offline policy evaluation for reinforcement learning. *In AISTATS'2021*.

Yin, M., Bai, Y., & Wang, Y. X. (2021). "Near-Optimal Offline Reinforcement Learning via Double Variance Reduction." *arXiv preprint arXiv:2102.01748* (2021).

Next lecture

- Notes on finite horizon MDP
- Some ideas behind how to improve the dependence on H or $1/(1-\gamma)$.
- RL algorithms:
 - Temporal difference learning
 - TD-learning with function approximation