

# **The Sample Complexity (with a Generative Model)**

**Sham Kakade and Wen Sun**

**CS 6789: Foundations of Reinforcement Learning**

# Announcements

- **Reading assignments** (see website)
  - sign up for a chapter (**signup sheep will be up today**)
  - start the assignment only after the we approve the chapter.
  - requirements:
    - one page report that summarizes the chapter
    - check **all** mathematical steps in the chapter
- **Participation/effort Bonus**
  - we will give extra credit for participation (class, ED, etc)
  - extra credit for reading assignments, finding bugs, project...
- **The book will be updated often.**
  - Feedback/questions/finding typos appreciated!

Today:

# Today:

- Recap: computational complexity
  - Question: Given an MDP  $\mathcal{M} = (S, A, P, r, \gamma)$  can we **exactly compute**  $Q^*$  (or find  $\pi^*$ ) in polynomial time?

# Today:

- Recap: computational complexity
  - Question: Given an MDP  $\mathcal{M} = (S, A, P, r, \gamma)$  can we **exactly compute**  $Q^*$  (or find  $\pi^*$ ) in polynomial time?
- Today: **statistical complexity**
  - Question: Given only sampling access to an unknown MDP  $\mathcal{M} = (S, A, P, r, \gamma)$  how many **observed transitions do we need to estimate**  $Q^*$  (or find  $\pi^*$ )?
  - Two sampling models: episodic setting and generative models.

Recap

# Summary Table

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

- VI: poly time for **fixed  $\gamma$** , not strongly poly
- PI: poly and strongly-poly time for **fixed  $\gamma$**
- LP approach: poly and strongly-poly time  
(LP approach is only logarithmic in  $1/(1 - \gamma)$ )

Today



## Two natural models for learning in an unknown MDP

- **Episodic setting:**

- in every episode,  $s_0 \sim \mu$ .
- the learner acts for some finite number of steps and observes the trajectory.
- The state is then resets to  $s_0 \sim \mu$ .

starling  
dist.



# Two natural models for learning in an unknown MDP

- **Episodic setting:**
  - in every episode,  $s_0 \sim \mu$ .
  - the learner acts for some finite number of steps and observes the trajectory.
  - The state is then resets to  $s_0 \sim \mu$ .
- **Generative model setting:**
  - input:  $(s, a)$
  - output: a sample  $s' \sim P(\cdot | s, a)$  and  $r(s, a)$

# Two natural models for learning in an unknown MDP

- **Episodic setting:**
  - in every episode,  $s_0 \sim \mu$ .
  - the learner acts for some finite number of steps and observes the trajectory.
  - The state is then resets to  $s_0 \sim \mu$ .
- **Generative model setting:**
  - input:  $(s, a)$
  - output: a sample  $s' \sim P(\cdot | s, a)$  and  $r(s, a)$
- **Sample complexity of RL:**  
how many transitions do we need observe in order to find a near optimal policy?
  - Episodic setting: we must actively explore to gather information
  - Generative model setting: lets us disentangle the issue of fundamental statistical limits from exploration.

# How many samples do we need to learn?

- What is the minmax optimal sample complexity, with generative modeling access?  
(using *any* algorithm)
  - Since  $P$  has  $S^2A$  parameters, we may hope that  $O(S^2A)$  samples are sufficient for learning.

# How many samples do we need to learn?

- What is the minmax optimal sample complexity, with generative modeling access? (using *any* algorithm)
  - Since  $P$  has  $S^2A$  parameters, we may hope that  $O(S^2A)$  samples are sufficient for learning.
- Questions:
  - Is a naive model-based approach optimal?  
i.e. estimate  $P$  accurately (using  $O(S^2A)$  samples) and then use  $\widehat{P}$  for planning.
  - Is sublinear learning possible?  
(i.e. learn with fewer than  $\Omega(S^2A)$  samples)

# How many samples do we need to learn?

- What is the minmax optimal sample complexity, with generative modeling access? (using *any* algorithm)
  - Since  $P$  has  $S^2A$  parameters, we may hope that  $O(S^2A)$  samples are sufficient for learning.
- Questions:
  - Is a naive model-based approach optimal?  
i.e. estimate  $P$  accurately (using  $O(S^2A)$  samples) and then use  $\widehat{P}$  for planning.
  - Is sublinear learning possible?  
(i.e. learn with fewer than  $\Omega(S^2A)$  samples)
- If sublinear learning is possible, then we do not need an accurate model of the world in order to act near-optimally?

# The most naive approach: model based

- Today: let us assume access to a [generative model](#)

# The most naive approach: model based

- Today: let us assume access to a **generative model**
- most naive approach to learning:
  - Call our simulator **N times at each state action pair.**
  - Let  $\widehat{P}$  be our empirical model:

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

where  $\text{count}(s', s, a)$  is the #times  $(s, a)$  transitions to state  $s'$ .

- we also know the rewards after one call.

(for simplicity, we often assume  $r(s, a)$  is deterministic)



# The most naive approach: model based

- Today: let us assume access to a **generative model**
- most naive approach to learning:
  - Call our simulator **N times at each state action pair.**

- Let  $\hat{P}$  be our empirical model:

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

$$\hat{P}(\cdot | s, a)$$

where  $\text{count}(s', s, a)$  is the #times  $(s, a)$  transitions to state  $s'$ .

- we also know the rewards after one call.

(for simplicity, we often assume  $r(s, a)$  is deterministic)

- The total number of calls to our generative model is **SAN**.

call the empirical MDP  $\hat{M}$

# Attempt 1:

the naive model based approach

# Model accuracy

**Proposition:**  $c$  is an absolute constant.  $\epsilon > 0$ . For  $N \geq \frac{c\gamma}{(1-\gamma)^4} \frac{S \log(cSA/\delta)}{\epsilon^2}$

and with probability greater than  $1 - \delta$ ,

# Model accuracy

$$\frac{c\gamma}{S^2 A}$$

**Proposition:**  $c$  is an absolute constant.  $\epsilon > 0$ . For  $N \geq \frac{c\gamma}{(1-\gamma)^4} \frac{S \log(cSA/\delta)}{\epsilon^2}$

and with probability greater than  $1 - \delta$ ,

$$\text{horizon} \approx \frac{1}{1-\gamma}$$

- Model accuracy: The transition model is  $\epsilon$  has error bounded as:

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

# Model accuracy

**Proposition:**  $c$  is an absolute constant.  $\epsilon > 0$ . For  $N \geq \frac{c\gamma}{(1-\gamma)^4} \frac{S \log(cSA/\delta)}{\epsilon^2}$

and with probability greater than  $1 - \delta$ ,

- Model accuracy: The transition model is  $\epsilon$  has error bounded as:

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

- Uniform value accuracy: For all policies  $\pi$ ,

$$\|Q^\pi - \widehat{Q}^\pi\|_\infty \leq \epsilon / 2$$

empirical  
MDP  $\mathcal{M}$   
 $\widehat{Q}^\pi$  is value of  $\pi$   
in  $\widehat{\mathcal{M}}$

# Model accuracy

**Proposition:**  $c$  is an absolute constant.  $\epsilon > 0$ . For  $N \geq \frac{c\gamma}{(1-\gamma)^4} \frac{S \log(cSA/\delta)}{\epsilon^2}$

and with probability greater than  $1 - \delta$ ,

- Model accuracy: The transition model is  $\epsilon$  has error bounded as:

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

- Uniform value accuracy: For all policies  $\pi$ ,

$$\|Q^\pi - \widehat{Q}^\pi\|_\infty \leq \epsilon / 2$$

- Near optimal planning: Suppose that  $\widehat{\pi}$  is the optimal policy in  $\widehat{M}$ .

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

size  $SA \times SA$  Matrix Expressions



- Define  $P^\pi$  to be the transition matrix on state-action pairs (for deterministic  $\pi$ ):

$$P_{(s,a),(s',a')}^\pi := \begin{cases} P(s' | s, a) & \text{if } a' = \pi(s') \\ 0 & \text{if } a' \neq \pi(s') \end{cases}$$

# Matrix Expressions

- Define  $P^\pi$  to be the transition matrix on state-action pairs (for deterministic  $\pi$ ):

$$P_{(s,a),(s',a')}^\pi := P(s' | s, a)$$

$$0$$

$$\text{if } a' = \pi(s')$$

$$\text{if } a' \neq \pi(s')$$

$P$  is of size  
 $SA \times S$

- With this notation,

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

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

$$Q^\pi \in \mathcal{R}^{SA}$$

$$V^\pi \in \mathcal{R}^S$$



# Matrix Expressions

- Define  $P^\pi$  to be the transition matrix on state-action pairs (for deterministic  $\pi$ ):

$$P_{(s,a),(s',a')}^\pi := \begin{cases} P(s' | s, a) & \text{if } a' = \pi(s') \\ 0 & \text{if } a' \neq \pi(s') \end{cases}$$

- With this notation,

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

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



- Also,

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

(where one can show the inverse exists)

$$r \in \mathbb{R}^{S \times A}$$

$$Q^\pi(s, a) = r(s, a)$$

$$+ \gamma \sum_{s' \sim P_{s, a}} [V^\pi(s')] ]$$

# “Simulation” Lemma

“Simulation Lemma”: For all  $\pi$ ,

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

# “Simulation” Lemma

“Simulation Lemma”: For all  $\pi$ ,

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

$$Q^\pi(s, \pi(s)) = V^\pi(s)$$

**Proof:** Using our matrix equality for  $Q^\pi$ , we have:

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

$$P^\pi Q^\pi = P V^\pi$$

# Proof of Claim 1

# Proof of Claim 1

- Concentration of a distribution in the  $\ell_1$  norm:
  - For a fixed  $s, a$ . With pr greater than  $1 - \delta$ ,

$$\|P(\cdot | s, a) - \widehat{P}(\cdot | s, a)\|_1 \leq c \sqrt{\frac{S \log(1/\delta)}{N}}$$

with  $N$  samples used to estimate  $\widehat{P}(\cdot | s, a)$ .

# Proof of Claim 1

- Concentration of a distribution in the  $\ell_1$  norm:
  - For a fixed  $s, a$ . With pr greater than  $1 - \delta$ ,

$$\|P(\cdot | s, a) - \widehat{P}(\cdot | s, a)\|_1 \leq c \sqrt{\frac{S \log(1/\delta)}{N}}$$

with  $N$  samples used to estimate  $\widehat{P}(\cdot | s, a)$ .

- The first claim now follows by the union bound.

over all SA dist.

$P(\cdot | s, a)$

# Proof of Claim 2 (&3)

## Proof of Claim 2 (&3)

For the second claim,

$$\begin{aligned}
 \|Q^\pi - \widehat{Q}^\pi\|_\infty &= \|\gamma(I - \gamma\widehat{P}^\pi)^{-1}(P - \widehat{P})V^\pi\|_\infty && \frac{1}{1-\gamma} (I - \gamma P^\pi)^{-1} \\
 &\leq \frac{\gamma}{1-\gamma} \|(P - \widehat{P})V^\pi\|_\infty && \max_{s,a} \left| \left[ P_{sa} V^\pi - \widehat{P}_{sa} V^\pi \right] \right| \\
 &\leq \frac{\gamma}{1-\gamma} \left( \max_{s,a} \|P(\cdot | s, a) - \widehat{P}(\cdot | s, a)\|_1 \right) \|V^\pi\|_\infty \\
 &\leq \frac{\gamma}{(1-\gamma)^2} \max_{s,a} \|P(\cdot | s, a) - \widehat{P}(\cdot | s, a)\|_1
 \end{aligned}$$

(why is the first inequality true?)



## Proof of Claim 2 (&3)

For the second claim,

$$\begin{aligned}\|Q^\pi - \widehat{Q}^\pi\|_\infty &= \|\gamma(I - \gamma\widehat{P}^\pi)^{-1}(P - \widehat{P})V^\pi\|_\infty \\ &\leq \frac{\gamma}{1 - \gamma} \|(P - \widehat{P})V^\pi\|_\infty \\ &\leq \frac{\gamma}{1 - \gamma} \left( \max_{s,a} \|P(\cdot | s, a) - \widehat{P}(\cdot | s, a)\|_1 \right) \|V^\pi\|_\infty \\ &\leq \frac{\gamma}{(1 - \gamma)^2} \max_{s,a} \|P(\cdot | s, a) - \widehat{P}(\cdot | s, a)\|_1\end{aligned}$$

(why is the first inequality true?)

The proof for the Claim 3 immediately follows from the second claim.

## Attempt 2:

obtaining sublinear sample complexity

idea: use concentration only on  $V^*$

$$\Omega(s^2 A)$$

# Reference sheet (defs/notation)

# Reference sheet (defs/notation)

- Remember: # samples from generative model =  $SAN$

# Reference sheet (defs/notation)

- Remember: # samples from generative model =  $SAN$
- $P^\pi$  is the transition matrix on state-action pairs for a deterministic policy  $\pi$ :  
$$P_{(s,a),(s',a')}^\pi := P(s' | s, a) \quad \text{if } a' = \pi(s')$$
$$0 \quad \text{if } a' \neq \pi(s')$$

# Reference sheet (defs/notation)

- Remember: # samples from generative model =  $SAN$
- $P^\pi$  is the transition matrix on state-action pairs for a deterministic policy  $\pi$ :

$$P_{(s,a),(s',a')}^\pi := \begin{cases} P(s' | s, a) & \text{if } a' = \pi(s') \\ 0 & \text{if } a' \neq \pi(s') \end{cases}$$

- With this notation,

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

# Reference sheet (defs/notation)

- Remember: # samples from generative model =  $SAN$

- $P^\pi$  is the transition matrix on state-action pairs for a deterministic policy  $\pi$ :

$$P_{(s,a),(s',a')}^\pi := P(s' | s, a) \quad \text{if } a' = \pi(s')$$

$$0 \quad \text{if } a' \neq \pi(s')$$

- With this notation,

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

- $\frac{1}{1-\gamma}(I - \gamma P^\pi)^{-1}$  is a matrix whose rows are probability distributions (why?)

# Reference sheet (defs/notation)

- Remember: # samples from generative model =  $SAN$

- $P^\pi$  is the transition matrix on state-action pairs for a deterministic policy  $\pi$ :

$$P_{(s,a),(s',a')}^\pi := P(s' | s, a) \quad \text{if } a' = \pi(s')$$

$$0 \quad \text{if } a' \neq \pi(s')$$

- With this notation,

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

- $\frac{1}{1 - \gamma} (I - \gamma P^\pi)^{-1}$  is a matrix whose rows are probability distributions (why?)

- $\widehat{Q}^*$ : optimal value in estimated model  $\widehat{M}$ .

$\widehat{\pi}^*$ : optimal policy in  $\widehat{M}$ .

$Q^{\widehat{\pi}^*}$ : (true) value of estimated policy.



# Attempt 2: Sublinear Sample Complexity

# Attempt 2: Sublinear Sample Complexity

**Proposition:** (Crude Value Bound) With probability greater than  $1 - \delta$ ,

$$\|Q^* - \widehat{Q}^*\|_\infty \leq \frac{\gamma}{(1-\gamma)^2} \sqrt{\frac{2 \log(2SA/\delta)}{N}}$$

$$\|Q^* - \widehat{Q}^{\pi^*}\|_\infty \leq \frac{\gamma}{(1-\gamma)^2} \sqrt{\frac{2 \log(2SA/\delta)}{N}}$$

*if  $N \approx O(\frac{1}{\epsilon^2})$*

*↘*

*$\approx O(\epsilon)$*

*$\Rightarrow \left(\frac{S}{\epsilon^2}\right)_{\text{samples}}^{\text{gen.}}$*

# Attempt 2: Sublinear Sample Complexity

**Proposition:** (Crude Value Bound) With probability greater than  $1 - \delta$ ,

$$\|Q^* - \widehat{Q}^*\|_\infty \leq \frac{\gamma}{(1-\gamma)^2} \sqrt{\frac{2 \log(2SA/\delta)}{N}}$$

$$\|Q^* - \widehat{Q}^{\pi^*}\|_\infty \leq \frac{\gamma}{(1-\gamma)^2} \sqrt{\frac{2 \log(2SA/\delta)}{N}}$$

*with amplification of  $\frac{1}{1-\gamma}$*

What about the value of the policy?

$$\|Q^* - Q^{\widehat{\pi}^*}\|_\infty \leq \frac{\gamma}{(1-\gamma)^3} \sqrt{\frac{2 \log(2SA/\delta)}{N}}$$

# Component-wise Bounds Lemma

**Lemma:** we have that

$$Q^* - \widehat{Q}^* \leq \gamma(I - \gamma \widehat{P}^{\pi^*})^{-1}(P - \widehat{P})V^*$$

$$Q^* - \widehat{Q}^* \geq \gamma(I - \gamma \widehat{P}^{\widehat{\pi}^*})^{-1}(P - \widehat{P})V^*$$

# Component-wise Bounds Lemma

Simulation Lemma:  $(P - \hat{P})V^*$

interpretation of  $(P - \hat{P})V^*$

**Lemma:** we have that

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

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

$$Q^* \wedge_{\pi^*} \hat{Q}^* = \hat{Q}^* \geq \hat{Q}^* \wedge_{\hat{\pi}^*} Q^*$$

**Proof:**

For the first claim, the optimality of  $\pi^*$  in  $M$  implies:

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

using the simulation lemma in the final step.

See notes for the proof of second claim.

Proof: (& key idea for sublinearity)

# Proof: (& key idea for sublinearity)

- Proof of the first claim:

# Proof: (& key idea for sublinearity)

- Proof of the first claim:

- By the previous lemma:  $\|Q^\star - \widehat{Q}^\star\|_\infty \leq \frac{\gamma}{1-\gamma} \|(P - \widehat{P})V^\star\|_\infty$



# Proof: (& key idea for sublinearity)

- Proof of the first claim:

- By the previous lemma:  $\|Q^\star - \widehat{Q}^\star\|_\infty \leq \frac{\gamma}{1-\gamma} \|(P - \widehat{P})V^\star\|_\infty$

- Recall  $\|V^\star\|_\infty \leq 1/(1-\gamma)$ .

# Proof: (& key idea for sublinearity)

- Proof of the first claim:

- By the previous lemma:  $\|Q^\star - \widehat{Q}^\star\|_\infty \leq \frac{\gamma}{1 - \gamma} \|(P - \widehat{P})V^\star\|_\infty$

- Recall  $\|V^\star\|_\infty \leq 1/(1 - \gamma)$ .

- By Hoeffding's inequality and the union bound,

$$\begin{aligned} \|(P - \widehat{P})V^\star\|_\infty &= \max_{s,a} \left| E_{s' \sim P(\cdot|s,a)}[V^\star(s')] - E_{s' \sim \widehat{P}(\cdot|s,a)}[V^\star(s')] \right| \\ &\leq \frac{1}{1 - \gamma} \sqrt{\frac{2 \log(2SA/\delta)}{N}} \end{aligned}$$

which holds with probability greater than  $1 - \delta$ .

# Proof: (& key idea for sublinearity)

- Proof of the first claim:

- By the previous lemma:  $\|Q^\star - \widehat{Q}^\star\|_\infty \leq \frac{\gamma}{1-\gamma} \|(P - \widehat{P})V^\star\|_\infty$

- Recall  $\|V^\star\|_\infty \leq 1/(1-\gamma)$ .

- By Hoeffding's inequality and the union bound,

$$\begin{aligned} \|(P - \widehat{P})V^\star\|_\infty &= \max_{s,a} \left| E_{s' \sim P(\cdot|s,a)}[V^\star(s')] - E_{s' \sim \widehat{P}(\cdot|s,a)}[V^\star(s')] \right| \\ &\leq \frac{1}{1-\gamma} \sqrt{\frac{2 \log(2SA/\delta)}{N}} \end{aligned}$$

which holds with probability greater than  $1 - \delta$ .

- Proof of second claim is similar (see the book)

# Attempt 3:

minimax optimal sample complexity

idea: better variance control

# (“near”) Minimax Optimal Sample Complexity

**Theorem:** (Azar et al. '13) With probability greater than  $1 - \delta$ ,

$$\|Q^* - \widehat{Q}^*\|_\infty \leq \gamma \sqrt{\frac{c}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{N}} + \frac{c\gamma}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{N},$$

where  $c$  is an absolute constant.

# (“near”) Minimax Optimal Sample Complexity

**Theorem:** (Azar et al. '13) With probability greater than  $1 - \delta$ ,

$$\|Q^* - \widehat{Q}^*\|_\infty \leq \gamma \sqrt{\frac{c}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{N}} + \frac{c\gamma}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{N},$$

where  $c$  is an absolute constant.

**Corollary:** for  $\epsilon < 1$ , provided  $N \geq \frac{c}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{\epsilon^2}$  then

$$\|Q^* - \widehat{Q}^*\|_\infty \leq \epsilon \text{ (with prob. greater than } 1 - \delta)$$

# (“near”) Minimax Optimal Sample Complexity

**Theorem:** (Azar et al. '13) With probability greater than  $1 - \delta$ ,

$$\|Q^* - \widehat{Q}^*\|_\infty \leq \gamma \sqrt{\frac{c}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{N}} + \frac{c\gamma}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{N},$$

where  $c$  is an absolute constant.

**Corollary:** for  $\epsilon < 1$ , provided  $N \geq \frac{c}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{\epsilon^2}$  then

$$\|Q^* - \widehat{Q}^*\|_\infty \leq \epsilon \text{ (with prob. greater than } 1 - \delta)$$

**Corollary:** What about the policy? Naively, need  $N/(1-\gamma)^2$  more samples. We pay *another factor of*  $1/(1-\gamma)^2$  samples. *Is this real?*

# Minimax Optimal Sample Complexity (on the policy)



# Minimax Optimal Sample Complexity (on the policy)

**Theorem:** (Agarwal et al. '20) For  $\epsilon < \sqrt{1/(1-\gamma)}$ , provided

$N \geq \frac{c}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{\epsilon^2}$  then with prob. greater than  $1 - \delta$ ,

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

# Minimax Optimal Sample Complexity (on the policy)

**Theorem:** (Agarwal et al. '20) For  $\epsilon < \sqrt{1/(1-\gamma)}$ , provided

$N \geq \frac{c}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{\epsilon^2}$  then with prob. greater than  $1 - \delta$ ,

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

**Lower Bound:** We can't do better.

# Proof sketch: part 1

- From “Component-wise Bounds” lemma, we want to bound:

$$Q^* - \widehat{Q}^* \leq \gamma \|(I - \gamma \widehat{P}^{\pi^*})^{-1} (P - \widehat{P}) V^*\|_{\infty} \leq ??$$

# Proof sketch: part 1

- From “Component-wise Bounds” lemma, we want to bound:

$$Q^* - \widehat{Q}^* \leq \gamma \|(I - \gamma \widehat{P}^{\pi^*})^{-1} (P - \widehat{P}) V^*\|_\infty \leq ??$$

- From Bernstein's ineq, with pr. greater than  $1 - \delta$ , we have (component-wise):

$$|(P - \widehat{P}) V^*| \leq \sqrt{\frac{2 \log(2SA/\delta)}{N}} \sqrt{\text{Var}_P(V^*)} + \frac{1}{1 - \gamma} \frac{2 \log(2SA/\delta)}{3N} \mathbf{1}$$

# Proof sketch: part 1

- From “Component-wise Bounds” lemma, we want to bound:

$$Q^* - \widehat{Q}^* \leq \gamma \|(I - \gamma \widehat{P}^{\pi^*})^{-1} (P - \widehat{P}) V^*\|_\infty \leq ??$$

- From Bernstein's ineq, with pr. greater than  $1 - \delta$ , we have (component-wise):

$$|(P - \widehat{P}) V^*| \leq \sqrt{\frac{2 \log(2SA/\delta)}{N}} \sqrt{\text{Var}_P(V^*)} + \frac{1}{1 - \gamma} \frac{2 \log(2SA/\delta)}{3N} \xrightarrow{1}$$

- Therefore

$$Q^* - \widehat{Q}^* \leq \gamma \sqrt{\frac{2 \log(2SA/\delta)}{N}} \|(I - \gamma \widehat{P}^{\pi^*})^{-1} \sqrt{\text{Var}_P(V^*)}\|_\infty \\ + \text{"lower order term"}$$

# Bellman Equation for the Variance

- **Variance:**  $\text{Var}_P(V)(s, a) := \text{Var}_{P(\cdot|s,a)}(V)$

Component wise variance:  $\text{Var}_P(V) := P(V)^2 - (PV)^2$

# Bellman Equation for the Variance

- **Variance:**  $\text{Var}_P(V)(s, a) := \text{Var}_{P(\cdot|s,a)}(V)$

Component wise variance:  $\text{Var}_P(V) := P(V)^2 - (PV)^2$

- Let's keep around the MDP M subscripts.

Define  $\Sigma_M^\pi$  as the (total) variance of the discounted reward:

$$\Sigma_M^\pi(s, a) := E \left[ \left( \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) - Q_M^\pi(s, a) \right)^2 \middle| s_0 = s, a_0 = a \right]$$

# Bellman Equation for the Variance

- **Variance:**  $\text{Var}_P(V)(s, a) := \text{Var}_{P(\cdot|s,a)}(V)$

Component wise variance:  $\text{Var}_P(V) := P(V)^2 - (PV)^2$

- Let's keep around the MDP M subscripts.

Define  $\Sigma_M^\pi$  as the (total) variance of the discounted reward:

$$\Sigma_M^\pi(s, a) := E \left[ \left( \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) - Q_M^\pi(s, a) \right)^2 \middle| s_0 = s, a_0 = a \right]$$

- **Bellman equation for the total variance:**

$$\Sigma_M^\pi = \gamma^2 \text{Var}_P(V_M^\pi) + \gamma^2 P^\pi \Sigma_M^\pi$$



# Key Lemma

**Lemma:** For any policy  $\pi$  and MDP  $M$ ,

$$\left\| (I - \gamma P^\pi)^{-1} \sqrt{\text{Var}_P(V_M^\pi)} \right\|_\infty \leq \sqrt{\frac{2}{(1 - \gamma)^3}}$$

Proof idea: convexity + Bellman equations for the variance.

# Putting it all together

**Proof sketch:** we have two MDPs  $M$  and  $\hat{M}$ . need to bound:

# Putting it all together

**Proof sketch:** we have two MDPs  $M$  and  $\widehat{M}$ . need to bound:

$$\|(I - \gamma \widehat{P}^{\pi^*})^{-1} \sqrt{\text{Var}_P(V^*)}\|_\infty = \|(I - \gamma P_{\widehat{M}}^{\pi^*})^{-1} \sqrt{\text{Var}_P(V_M^{\pi^*})}\|_\infty$$

$$\leq \|(I - \gamma P_{\widehat{M}}^{\pi^*})^{-1} \sqrt{\text{Var}_P(V_{\widehat{M}}^{\pi^*})}\| + \text{"lower order"}$$

$$\leq \sqrt{\frac{2}{(1 - \gamma)^3}} + \text{"lower order"}$$

# Putting it all together

**Proof sketch:** we have two MDPs  $M$  and  $\widehat{M}$ . need to bound:

$$\|(I - \gamma \widehat{P}^{\pi^*})^{-1} \sqrt{\text{Var}_P(V^*)}\|_\infty = \|(I - \gamma P_{\widehat{M}}^{\pi^*})^{-1} \sqrt{\text{Var}_P(V_{\widehat{M}}^{\pi^*})}\|_\infty$$

$$\leq \|(I - \gamma P_{\widehat{M}}^{\pi^*})^{-1} \sqrt{\text{Var}_P(V_{\widehat{M}}^{\pi^*})}\| + \text{"lower order"}$$

$$\leq \sqrt{\frac{2}{(1 - \gamma)^3}} + \text{"lower order"}$$

First equality above: just notation

# Putting it all together

**Proof sketch:** we have two MDPs  $M$  and  $\widehat{M}$ . need to bound:

$$\|(I - \gamma \widehat{P}^{\pi^*})^{-1} \sqrt{\text{Var}_P(V^*)}\|_\infty = \|(I - \gamma P_{\widehat{M}}^{\pi^*})^{-1} \sqrt{\text{Var}_P(V_M^{\pi^*})}\|_\infty$$

$$\leq \|(I - \gamma P_{\widehat{M}}^{\pi^*})^{-1} \sqrt{\text{Var}_P(V_{\widehat{M}}^{\pi^*})}\| + \text{"lower order"}$$

$$\leq \sqrt{\frac{2}{(1 - \gamma)^3}} + \text{"lower order"}$$

First equality above: just notation

Second step: concentration  $\rightarrow$  we need to quantify:

$$\sqrt{\text{Var}_P(V_M^{\pi^*})} \approx \sqrt{\text{Var}_P(V_{\widehat{M}}^{\pi^*})}$$

# Putting it all together

**Proof sketch:** we have two MDPs  $M$  and  $\widehat{M}$ . need to bound:

$$\|(I - \gamma \widehat{P}^{\pi^*})^{-1} \sqrt{\text{Var}_P(V^*)}\|_\infty = \|(I - \gamma P_{\widehat{M}}^{\pi^*})^{-1} \sqrt{\text{Var}_P(V_M^{\pi^*})}\|_\infty$$

$$\leq \|(I - \gamma P_{\widehat{M}}^{\pi^*})^{-1} \sqrt{\text{Var}_P(V_{\widehat{M}}^{\pi^*})}\| + \text{"lower order"}$$

$$\leq \sqrt{\frac{2}{(1 - \gamma)^3}} + \text{"lower order"}$$

First equality above: just notation

Second step: concentration  $\rightarrow$  we need to quantify:

$$\sqrt{\text{Var}_P(V_M^{\pi^*})} \approx \sqrt{\text{Var}_P(V_{\widehat{M}}^{\pi^*})}$$

Last step: previous slide