The Sample Complexity with a Generative Model

Announcements

- Norms for class
 - Video: please feel free to turn your video! It's nice to see some of you once in a while.
 - Questions: by chat or just ask.
- Notes are posted in advance. Notes/Book on hand can help during lectures.
- HW1 posted later this week (due in two weeks)

Today:

- Recap: computational complexity
 - Question: Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$ can we exactly compute Q^* (or find π^*) in polynomial time?
- Today: statistical complexity
 - Question: Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$ how many observed transitions do we need to estimate Q^* (or find π^*)?
 - We consider an abstract model (a generative model) to study the number of samples required for learning.

Recap

Summary Table

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

- VI Per iteration complexity: S^2A
- PI Per iteration complexity: $S^3 + S^2A$
- The LP approach is only logarithmic in $1-\gamma$

Today

• A generative model provides us with a sample $s' \sim P(\cdot | s, a)$ upon input of a state action pair (s, a). Know the reward function r(s, a)

- A generative model provides us with a sample s' ∼ P(· | s, a) upon input of a state action pair (s, a).
- most naive approach to learning: suppose we call our simulator N times at each state action pair. Let \widehat{P} be our empirical model:

- A generative model provides us with a sample s' ∼ P(· | s, a) upon input of a state action pair (s, a).
- most naive approach to learning: suppose we call our simulator N times at each state action pair. Let P be our empirical model:
 P (s' | s, a) = count(s', s, a) / N

where count(s', s, a) is the #times the state-action (s, a) transitions to state s'.

- A generative model provides us with a sample s' ∼ P(· | s, a) upon input of a state action pair (s, a).
- most naive approach to learning: suppose we call our simulator N times at each state action pair. Let P be our empirical model:
 P (s' | s, a) = count(s', s, a) / N

where count(s', s, a) is the #times the state-action (s, a) transitions to state s'.

• Each "sample" calls the generative model SA times (for each state action pair). The total number of calls to our generative model is SAN.

- A generative model provides us with a sample $s' \sim P(\cdot | s, a)$ upon input of a state action pair (s, a).
- most naive approach to learning: suppose we call our simulator N times at each state action pair. Let P be our empirical model:
 P (s' | s, a) = count(s', s, a) / N

where count(s', s, a) is the #times the state-action (s, a) transitions to state s'.

- Each "sample" calls the generative model *SA* times (for each state action pair). The total number of calls to our generative model is *SAN*.
- The generative model helps us disentangle the issue of fundamental statistical lin from exploration.



• This is a reasonable abstraction to understand the statistical limit, without having to directly address exploration.

How many samples do we need to learn?

- This is a reasonable abstraction to understand the statistical limit, without having to directly address exploration.
 Note that since *P* has a S²A parameters, a naive approach would be to
- Note that since *P* has a S^2A parameters, a naive approach would be to estimate *P* accurately (using $O(S^2A)$ samples) and then use \widehat{P} for planning.

How many samples do we need to learn?

- This is a reasonable abstraction to understand the statistical limit, without having to directly address exploration.
- Note that since *P* has a S^2A parameters, a naive approach would be to estimate *P* accurately (using $O(S^2A)$ samples) and then use \widehat{P} for planning.
- Do we require an accurate model of the world in order to find a near optimal policy?
 Scoptimal
 Model
 Scoptimal
 Model
 Model</l



Model accuracy # SAN times

0(54/22)

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

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

Model accuracy

Proposition: c is an absolute constant. $\epsilon > 0$. For $N \ge \frac{c\gamma}{(1-\gamma)^4} \frac{S \log(cSA/\delta)}{\epsilon^2}$ and with probability greater than $1 - \delta$, • Model accuracy: The transition model is ϵ has error bounded as: $\max_{s,a} \|P(\cdot \mid s, a) - \widehat{P}(\cdot \mid s, a)\|_1 \le (1-\gamma)^2 \epsilon/2.$



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

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

- Model accuracy: The transition model is ϵ has error bounded as: ٠ $\max_{s,a} \|P(\cdot \mid s, a) - \widehat{P}(\cdot \mid s, a)\|_1 \le (1 - \gamma)^2 \epsilon/2.$
- Uniform value accuracy: For all policies π , $\|Q^{\pi} - \widehat{Q}^{\pi}\|_{\infty} \le \epsilon/2$

Model accuracy (STA) samples.

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

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

- Model accuracy: The transition model is ϵ has error bounded as: $\max \|P(\cdot | s, a) - \widehat{P}(\cdot | s, a)\|_1 \le (1 - \gamma)^2 \epsilon/2.$ value of the policy we find. s.a
- Uniform value accuracy: For all policies π ,
- $$\begin{split} \|Q^{\pi} \widehat{Q}^{\pi}\|_{\infty} &\leq \epsilon/2 \\ \bullet \text{ Near optimal planning: Suppose that } \widehat{\pi} \text{ is the optimal policy in } \widehat{M}. \\ \|\widehat{Q}^{\pi} Q^{\star}\|_{\infty} &\leq \epsilon \quad |Q^{\star} \sim Q^{\star}|| \leq \epsilon \quad |Q^{\star} Q^{\star}|| \leq \epsilon \end{split}$$

Matrix Expressions

• Define P^{π} to be the transition matrix on state-action pairs (for deterministic π):

 $P^{\pi}_{(s,a),(s',a')} := \underbrace{P(s' \mid s, a)}_{0} \quad \text{if } a' = \pi(s') \qquad \qquad ? \stackrel{\text{fi}}{\longrightarrow} \quad 5 A \times S A$ P matrix SAXS

Matrix Expressions

• Define P^{π} to be the transition matrix on state-action pairs (for deterministic π):

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

$$0 \quad \text{if } a' \neq \pi(s')$$
With this notation,
$$\sqrt{e} p^{\sharp} \quad \sqrt{e} q^{\xi} \wedge \sqrt{e} q^{\xi}$$

• With this notation, $Q^{\pi} = r + PV^{\pi}$ $Q^{\pi} = r + P^{\pi}Q^{\pi}$

Matrix Expressions

• Define P^{π} to be the transition matrix on state-action pairs (for deterministic π):

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

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

• With this notation, $Q^{\pi} = r + PV^{\pi}$ $Q^{\pi} = r + P^{\pi}Q^{\pi}$

$$= \sum (I - \gamma P^{\tau}) Q^{\tau} = \sum$$

 $Q^{\overline{n}} = (I - r \hat{p}^{\overline{n}})'r$

• And that:

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

(where one can show the inverse exists)

"Simulation" Lemma

(Simulation Lemma) For all π , we have: $Q^{\pi} - \widehat{Q}^{\pi} = \gamma (I - \gamma \widehat{P}^{\pi})^{-1} (P - \widehat{P}) V^{\pi}$

"Simulation" Lemma

(Simulation Lemma) For all
$$\pi$$
, we have:
 $Q^{\pi} - \widehat{Q}^{\pi} = \gamma (I - \gamma \widehat{P}^{\pi})^{-1} (P - \widehat{P}) V^{\pi}$
proof: Using our matrix equality for Q^{π} , we have:
 $Q^{\pi} - \widehat{Q}^{\pi} = (I - \gamma P^{\pi})^{-1} r - (I - \gamma \widehat{P}^{\pi})^{-1} r = Q^{\pi} - (\mathcal{I} - \gamma \widehat{P}^{\pi})^{-1} (I - \gamma \widehat{P}^{\pi})^{-1} r = Q^{\pi} - (\mathcal{I} - \gamma \widehat{P}^{\pi})^{-1} (I - \gamma \widehat{P}^{\pi}) - (I - \gamma P^{\pi})) Q^{\pi}$
 $= \gamma (I - \gamma \widehat{P}^{\pi})^{-1} (P^{\pi} - \widehat{P}^{\pi}) Q^{\pi}$ when $m \not f$. \widehat{Q}^{π}
 $= \gamma (I - \gamma \widehat{P}^{\pi})^{-1} (P - \widehat{P}) V^{\pi}$

Proof: Claim 1

Concentration of a distribution in the
$$\ell_1$$
 norm: for a fixed s, a
 $\|P(\cdot | s, a) - \widehat{P}(\cdot | s, a)\|_1 \leq c \sqrt{\frac{S \log(1/\delta)}{m}}$
with pr greater than $1 - \delta$, with m samples used to estimate $\widehat{P}(\cdot | s, a)$.

The first claim now follows by the union bound.



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

Reference sheet (defs/notation)

- Remember: # samples from generative model = SAN
- $\begin{array}{ll} \cdot P^{\pi} \text{ is the transition matrix on state-action pairs for a deterministic policy } \pi:\\ P^{\pi}_{(s,a),(s',a')} := P(s' \mid s, a) & \text{ if } a' = \pi(s') \\ 0 & \text{ if } a' \neq \pi(s') \\ \end{array} \\ \begin{array}{ll} \cdot \text{With this notation,} \\ Q^{\pi}_{,} = r + PV^{\pi}, \quad Q^{\pi} = r + P^{\pi}Q^{\pi}, \quad Q^{\pi} = (I \gamma P^{\pi})^{-1}r \end{array}$

• $(I - \gamma P^{\pi})^{-1}$ is a matrix whose rows are probability distributions (why?) • Estimated transition \widehat{P} , optimal value in estimated model \widehat{Q}^{\star} , optimal policy in optimal model $\widehat{\pi}^{\star}$, (true) value of estimated policy $Q^{\widehat{\pi}^{\star}}$



$\begin{array}{c} \mathcal{OSA} \\ \text{Attempt 2: calls to estimate} \\ \text{sublinear sample complexity!} \end{array}$

Attempt 2: Sublinear Sample Complexity

Attempt 2: Sublinear Sample Complexity

Proposition: (Crude Value Bound) Let $\delta \ge 0$. With probability greater than $1 - \delta$,

Attempt 2: Sublinear Sample Complexity

Proposition: (Crude Value Bound) Let $\delta \geq 0$. With probability greater than $1 - \delta$,

a a + t a

scaling ff $Q^{*} \in O(\frac{f''}{1-Y})$

$$\begin{split} \|Q^{\star} - \widehat{Q}^{\star}\|_{\infty} &\leq \frac{\gamma}{(1-\gamma)^2} \sqrt{\frac{2\log(2SA/\delta)}{N}} \\ \|Q^{\star} - \widehat{Q}^{\pi^{\star}}\|_{\infty} &\leq \frac{\gamma}{(1-\gamma)^2} \sqrt{\frac{2\log(2SA/\delta)}{N}} \\ &\alpha n\rho \|f\|_{\infty} \leq \frac{\gamma}{(1-\gamma)^2} \sqrt{\frac{2\log(2SA/\delta)}{N}} \\ \end{split}$$
What about the value of the policy?
$$\|Q^{\star} - Q^{\widehat{\pi}^{\star}}\|_{\infty} \leq \frac{\gamma}{(1-\gamma)^3} \sqrt{\frac{2\log(2SA/\delta)}{N}} \end{split}$$

Component-wise Bounds Lemma $a \le b^2$

Lemma: we have that:

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

Component-wise Bounds Lemma

Lemma: we have that:

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

Component-wise Bounds Lemma

 $\left(P-\rho\right)$

620

 \rightarrow $\hat{h}^{\hat{h}^*} 2 \hat{o}^{\pi^*}$

Estappine (V*) - Elve

Lemma: we have that:

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

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

proof:

For the first claim, the optimality of π^* in *M* implies:

 $Q^{\star} - \widehat{Q}^{\star} = Q^{\pi^{\star}} - \widehat{Q}^{\pi^{\star}} \leq Q^{\pi^{\star}} - \widehat{Q}^{\pi^{\star}} = \gamma (I - \gamma \widehat{P}^{\pi^{\star}})^{-1} (P - \widehat{P}) V^{\star}.$ using the simulation lemma in the final step.

ta Off ZQA See notes for the proof of second claim.

Proof: (& key idea for sublinearity!)

For the first claim,

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

Why?

Proof: (& key idea for sublinearity!)

For the first claim.

For the first claim,
$$\|Q^{\star} - \widehat{Q}^{\star}\|_{\infty} \leq \frac{\gamma}{1 - \gamma} \|(P - \widehat{P})V^{\star}\|_{\infty} \qquad \stackrel{\frown}{\bigcirc} \leq \overrightarrow{\zeta} \leq \overbrace{1}^{\frown}$$

Why?

By applying Hoeffding's inequality V^{\star} and the union bound, $\|(P - \widehat{P})V^{\star}\|_{\infty} = \max_{s,a} \left| E_{s' \sim P(\cdot|s,a)}[V^{\star}(s')] - E_{s' \sim \widehat{P}(\cdot|s,a)}[V^{\star}(s')] \right|$ $\leq \frac{1}{1-\gamma} \sqrt{\frac{2\log(2SA/\delta)}{N}}$

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

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: What about the policy? Need $N/(1 - \gamma)^2$ 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 $N \ge \frac{c}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{\epsilon^2}$ then with prob. greater than $1-\delta$), $Q^{*} \approx O(\overline{D})$ $\|Q^{\star} - Q^{\hat{\pi}\star}\|_{\infty} \le \epsilon$ Suppose $||Q^* - Q^{\frac{1}{n}}||_{\infty} \leq \frac{2}{1-3}$ sufficient NZ (SA)

Minimax Optimal Sample Complexity (on the policy)

Theorem: (Agarwal et al. '20) For $\epsilon < \sqrt{1/(1-\gamma)}$, provided $N \ge \frac{c}{(1-\gamma)^3} \frac{\log(cSA/\delta)}{\epsilon^2}$ then with prob. greater than $1 - \delta$), $\|Q^{\star} - Q^{\hat{\pi}\star}\|_{\infty} \le \epsilon$

Lower Bound: We can't do better.

Proof sketch: part 1

• From "Component-wise Bounds" lemma, we want to bound: $\|Q^{\star} - \widehat{Q}^{\star}\|_{\infty} \leq \gamma \|(I - \gamma \widehat{P}^{\pi^{\star}})^{-1}(P - \widehat{P})V^{\star}\|_{\infty} \leq ??$

Proof sketch: part 1

- From "Component-wise Bounds" lemma, we want to bound: $\|Q^{\star} - \widehat{Q}^{\star}\|_{\infty} \leq \gamma \|(I - \gamma \widehat{P}^{\pi^{\star}})^{-1}(P - \widehat{P})V^{\star}\|_{\infty} \leq ??$
- From Bernstein's ineq, with pr. greater than 1δ , we have (component-wise): $|(P - \widehat{P})V^{\star}| \leq \sqrt{\frac{2\log(2SA/\delta)}{N}}\sqrt{\operatorname{Var}_{P}(V^{\star})} + \frac{1}{1 - \gamma}\frac{2\log(2SA/\delta)}{3N}\overrightarrow{1}$

Proof sketch: part 1

- From "Component-wise Bounds" lemma, we want to bound: $\|Q^{\star} - \widehat{Q}^{\star}\|_{\infty} \leq \gamma \|(I - \gamma \widehat{P}^{\pi^{\star}})^{-1}(P - \widehat{P})V^{\star}\|_{\infty} \leq ??$
- From Bernstein's ineq, with pr. greater than 1δ , we have (component-wise): $|(P - \widehat{P})V^{\star}| \leq \sqrt{\frac{2\log(2SA/\delta)}{N}}\sqrt{\operatorname{Var}_{P}(V^{\star})} + \frac{1}{1 - \gamma} \frac{2\log(2SA/\delta)}{3N} \overrightarrow{1}$
- Therefore

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

Bellman Equation for the Variance

• Variance: $\operatorname{Var}_P(V)(s, a) := \operatorname{Var}_{P(\cdot|s,a)}(V)$ Component wise variance: $\operatorname{Var}_P(V) := P(V)^2 - (PV)^2$

Bellman Equation for the Variance

• Variance: $\operatorname{Var}_P(V)(s, a) := \operatorname{Var}_{P(\cdot|s,a)}(V)$ Component wise variance: $\operatorname{Var}_P(V) := P(V)^2 - (PV)^2$

• Let's keep around the MDP M subscripts. Define Σ_M^{π} as the (totoal) 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: $\operatorname{Var}_P(V)(s, a) := \operatorname{Var}_{P(\cdot|s,a)}(V)$ Component wise variance: $\operatorname{Var}_P(V) := P(V)^2 - (PV)^2$

- Let's keep around the MDP M subscripts. Define Σ_M^{π} as the (totoal) 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 π and MDP M,

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

Proof idea: convexity + Bellman equations for the variance.

Putting it all together Proof: we have two MDPs M and \widehat{M} . need to bound:

Proof: we have two MDPs M and M. need to bound:

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

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

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

Proof: we have two MDPs M and M. need to bound:

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

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

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

First equality above: just notation

Proof: we have two MDPs M and M. need to bound:

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

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

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

First equality above: just notation Second step: concentration -> need to say

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

Proof: we have two MDPs M and M. need to bound:

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

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

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

First equality above: just notation Second step: concentration -> need to say

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

Last step: previous slide