

# **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  $\pi^*$ ) 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  $\pi^*$ )?
  - 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.	$S^2 A \frac{L(P,r,\gamma) \log \frac{1}{1-\gamma}}{1-\gamma}$	$(S^3 + S^2 A) \frac{L(P,r,\gamma) \log \frac{1}{1-\gamma}}{1-\gamma}$	$S^3 A L(P, r, \gamma)$
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 Per iteration complexity:  $S^2 A$
- PI Per iteration complexity:  $S^3 + S^2 A$
- The LP approach is only logarithmic in  $1 - \gamma$

Today

# A Generative Model

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

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

where  $\text{count}(s', s, a)$  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 learning from 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  $S^2A$  parameters, a naive approach would be to estimate  $P$  accurately (using  $O(S^2A)$  samples) and then use  $\hat{P}$  for planning.
- *Do we require an accurate model of the world in order to find a near optimal policy?*



**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: 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  $\widehat{M}$ .

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

# Matrix Expressions

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

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

- With this notation,

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

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

- And that:

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

(where one can show the inverse exists)

# “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  $Q^\pi$ , we have:

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

# 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.

# Proof: Claim 2

For the second claim, we have that:

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

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

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

(why is step 1 true?)

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

# Reference sheet (defs/notation)

- Remember: # samples from generative model =  $SAN$
- $P^\pi$  is the transition matrix on state-action pairs for a deterministic policy  $\pi$ :  
$$P^\pi_{(s,a),(s',a')} := P(s' | s, a) \quad \text{if } a' = \pi(s')$$
$$0 \quad \text{if } a' \neq \pi(s')$$
- With this notation,  
$$Q^\pi = r + PV^\pi, \quad Q^\pi = r + P^\pi Q^\pi, \quad Q^\pi = (I - \gamma P^\pi)^{-1} r$$
- $\frac{1}{1 - \gamma} (I - \gamma P^\pi)^{-1}$  is a matrix whose rows are probability distributions (why?)
- Estimated transition  $\hat{P}$ , optimal value in estimated model  $\hat{Q}^*$ ,  
optimal policy in optimal model  $\hat{\pi}^*$ , (true) value of estimated policy  $Q^{\hat{\pi}^*}$

Attempt 2:

sublinear sample complexity!

idea: use concentration only on  $V^\star$



# Attempt 2: Sublinear Sample Complexity

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

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

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

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:

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

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

proof:

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

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

using the simulation lemma in the final step.

See notes for the proof of second claim.

# Proof: (& key idea for sublinearity!)

For the first claim,

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

Why?

By applying Hoeffding's inequality  $V^\star$  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  $1 - \delta$ .

# 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.

**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? 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)

**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:

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

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

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

- Therefore

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

+ "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:  $\text{Var}_P(V) := P(V)^2 - (PV)^2$

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

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

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

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

First equality above: just notation

Second step: concentration -> need to say

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

Last step: previous slide