

# **Demystifying the efficiency of reinforcement learning: A few recent stories**

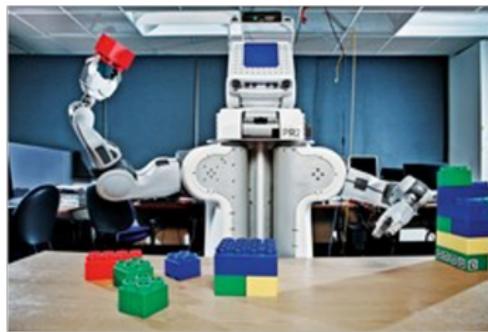


Yuxin Chen

EE, Princeton University

# Reinforcement learning (RL)

---



# RL challenges

---

In RL, an agent learns by interacting with an environment

- unknown or changing environments
- delayed rewards or feedback
- enormous state and action space
- nonconvexity



# Sample efficiency

---

Collecting data samples might be expensive or time-consuming



clinical trials



online ads

# Sample efficiency

---

Collecting data samples might be expensive or time-consuming



clinical trials



online ads

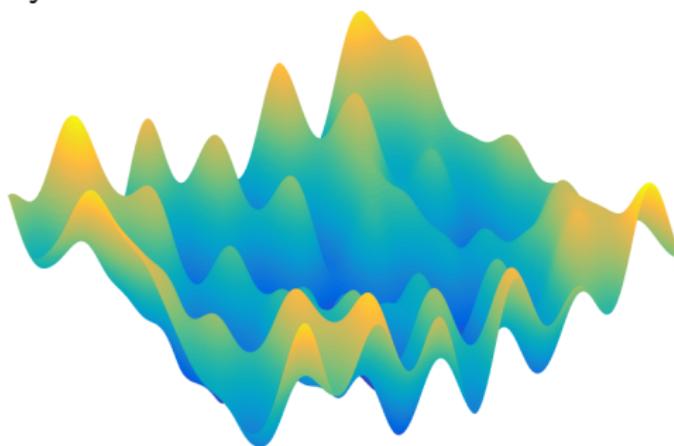
Calls for design of sample-efficient RL algorithms!

# Computational efficiency

---

Running RL algorithms might take a long time ...

- enormous state-action space
- nonconvexity

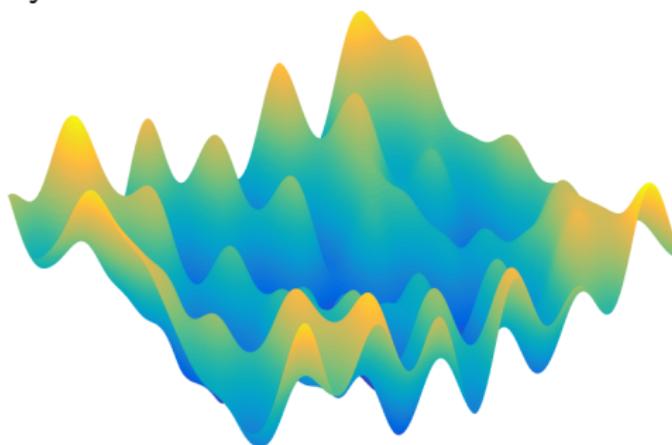


# Computational efficiency

---

Running RL algorithms might take a long time ...

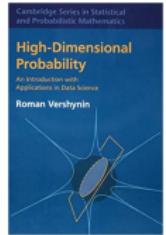
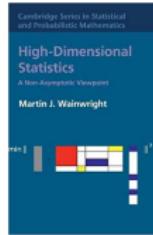
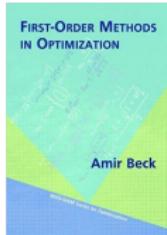
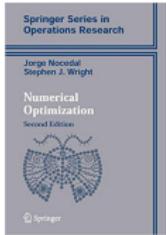
- enormous state-action space
- nonconvexity



Calls for computationally efficient RL algorithms!

# This talk: three recent stories

---



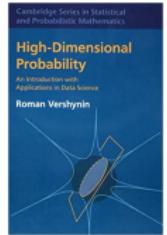
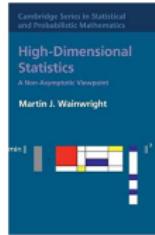
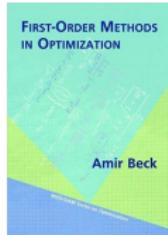
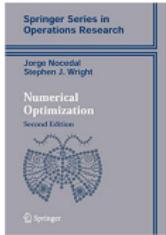
(large-scale) optimization

(high-dimensional) statistics

Demystify **sample-** and **computational** efficiency of RL algorithms

# This talk: three recent stories

---



(large-scale) optimization

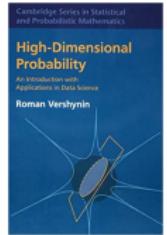
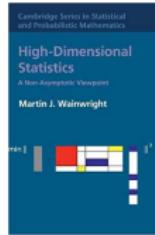
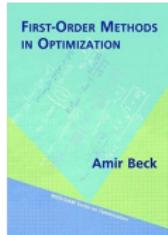
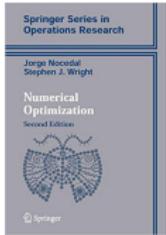
(high-dimensional) statistics

Demystify **sample-** and **computational** efficiency of RL algorithms

1. **model-based RL**
2. **policy-based RL**
3. **value-based RL**

# This talk: three recent stories

---



(large-scale) optimization

(high-dimensional) statistics

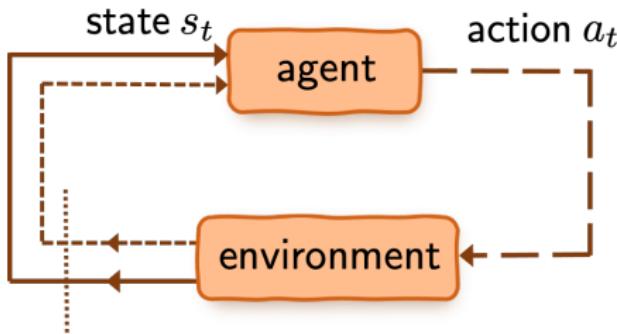
Demystify **sample-** and **computational** efficiency of RL algorithms

1. **model-based RL**: breaking a sample size barrier
2. **policy-based RL**: natural policy gradient (NPG) methods
3. **value-based RL**: Q-learning over Markovian samples

*Background: Markov decision processes*

# Markov decision process (MDP)

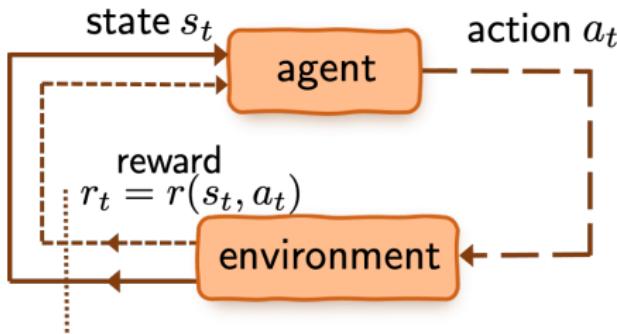
---



- $\mathcal{S}$ : state space
- $\mathcal{A}$ : action space

# Markov decision process (MDP)

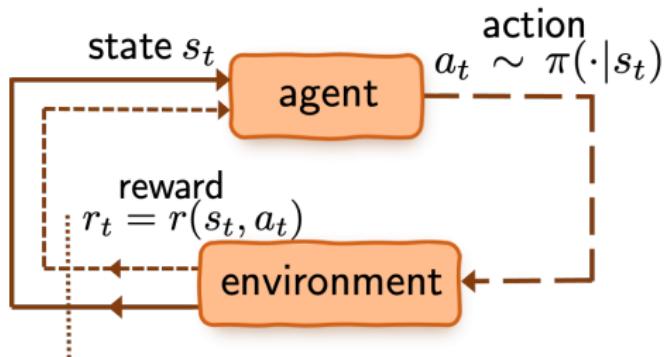
---



- $\mathcal{S}$ : state space
- $\mathcal{A}$ : action space
- $r(s, a) \in [0, 1]$ : immediate reward

# Markov decision process (MDP)

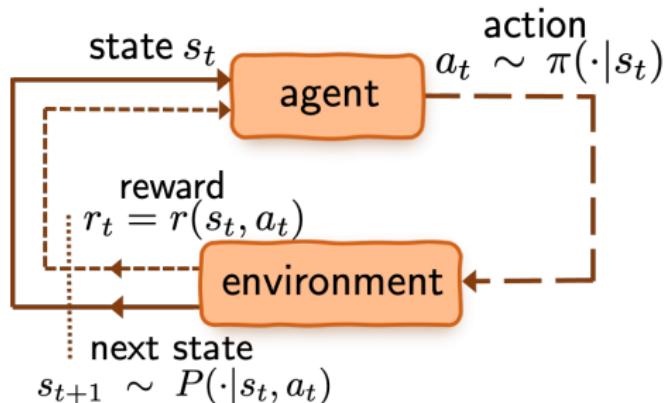
---



- $\mathcal{S}$ : state space
- $\mathcal{A}$ : action space
- $r(s, a) \in [0, 1]$ : immediate reward
- $\pi(\cdot|s)$ : policy (or action selection rule)

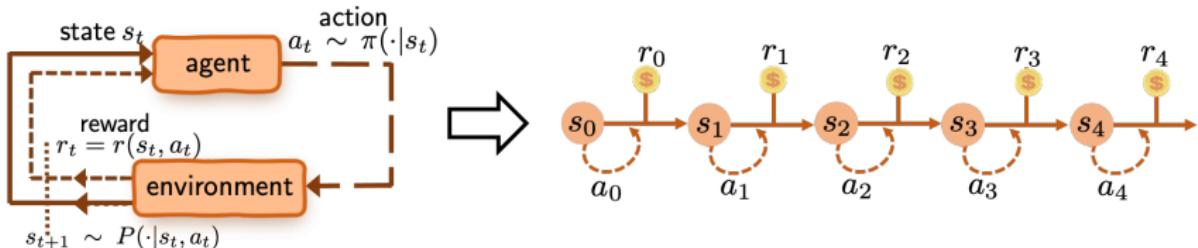
# Markov decision process (MDP)

---



- $\mathcal{S}$ : state space
- $\mathcal{A}$ : action space
- $r(s, a) \in [0, 1]$ : immediate reward
- $\pi(\cdot|s)$ : policy (or action selection rule)
- $P(\cdot|s, a)$ : **unknown** transition probabilities

# Value function

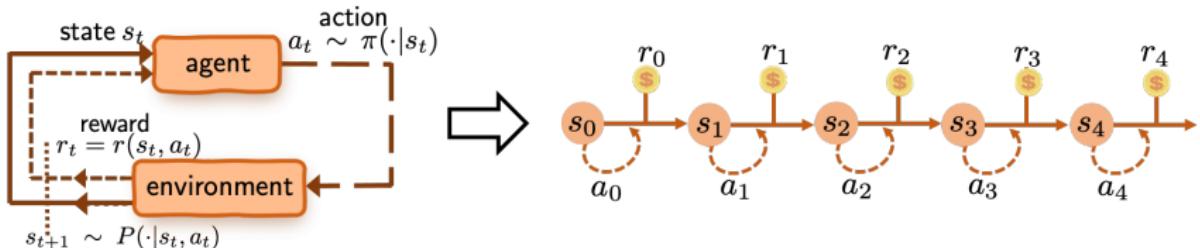


Value of policy  $\pi$ : long-term *discounted* reward

$$\forall s \in \mathcal{S} : V^\pi(s) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right]$$



# Value function



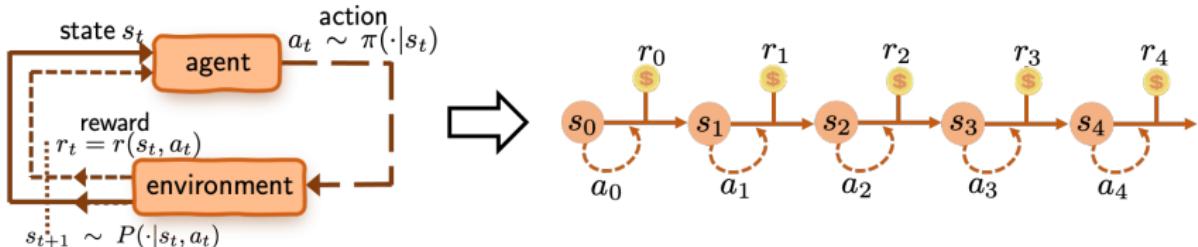
Value of policy  $\pi$ : long-term *discounted* reward

$$\forall s \in \mathcal{S} : V^\pi(s) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right]$$



- $(a_0, s_1, a_1, s_2, a_2, \dots)$ : generated under policy  $\pi$

# Value function



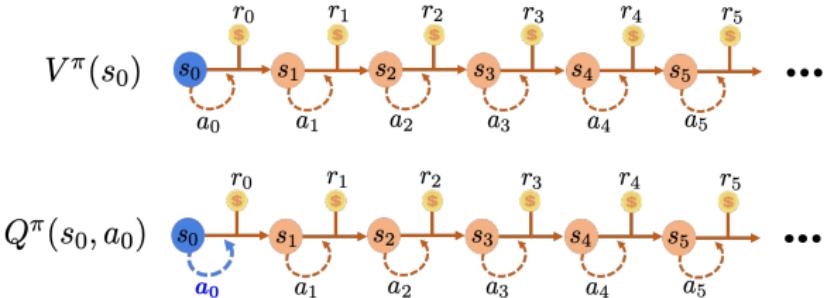
Value of policy  $\pi$ : long-term *discounted* reward

$$\forall s \in \mathcal{S} : V^\pi(s) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right]$$



- $(a_0, s_1, a_1, s_2, a_2, \dots)$ : generated under policy  $\pi$
- $\gamma \in [0, 1]$ : discount factor
  - take  $\gamma \rightarrow 1$  to approximate *long-horizon* MDPs

# Q-function



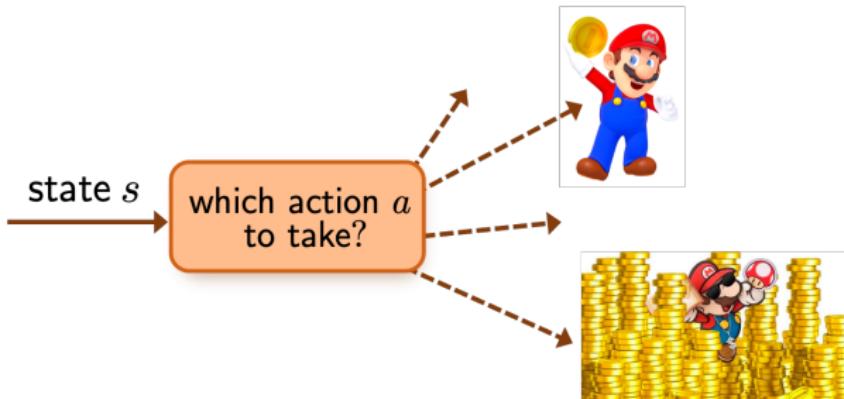
Q-function of policy  $\pi$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad Q^\pi(s, a) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

- $(\cancel{a_0}, s_1, a_1, s_2, a_2, \dots)$ : generated under policy  $\pi$

# Optimal policy and optimal values

---



- **Optimal policy  $\pi^*$ :** maximizing the value function

# Optimal policy and optimal values

---



- **Optimal policy**  $\pi^*$ : maximizing the value function
- Optimal values:  $V^* := V^{\pi^*}$

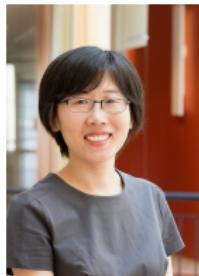
*Story 1: breaking the sample size barrier  
via **model-based RL** under a generative model*



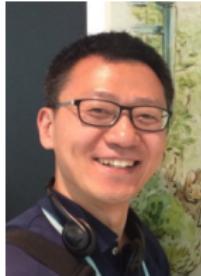
Gen Li  
Tsinghua EE



Yuting Wei  
CMU Stats

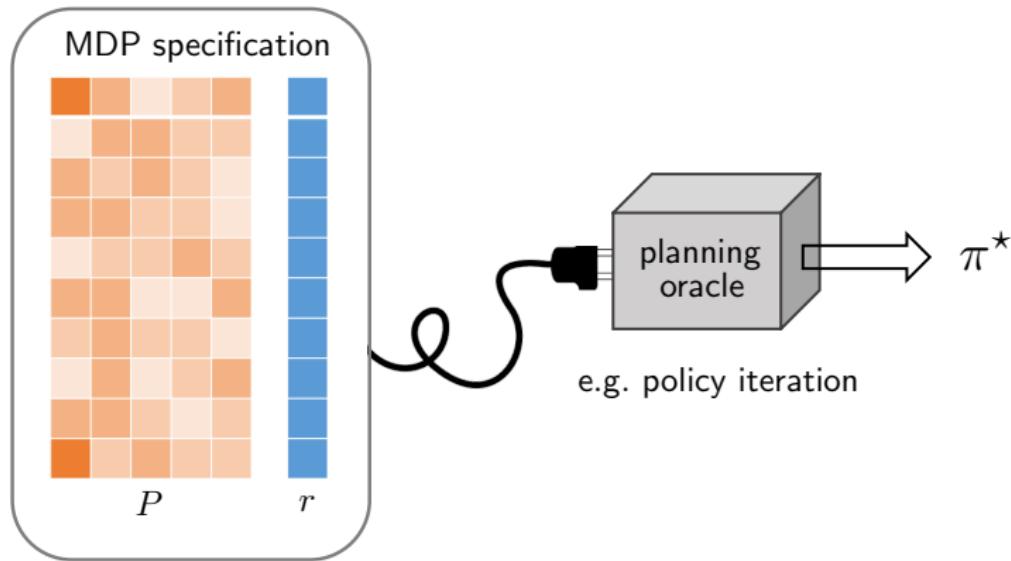


Yuejie Chi  
CMU ECE



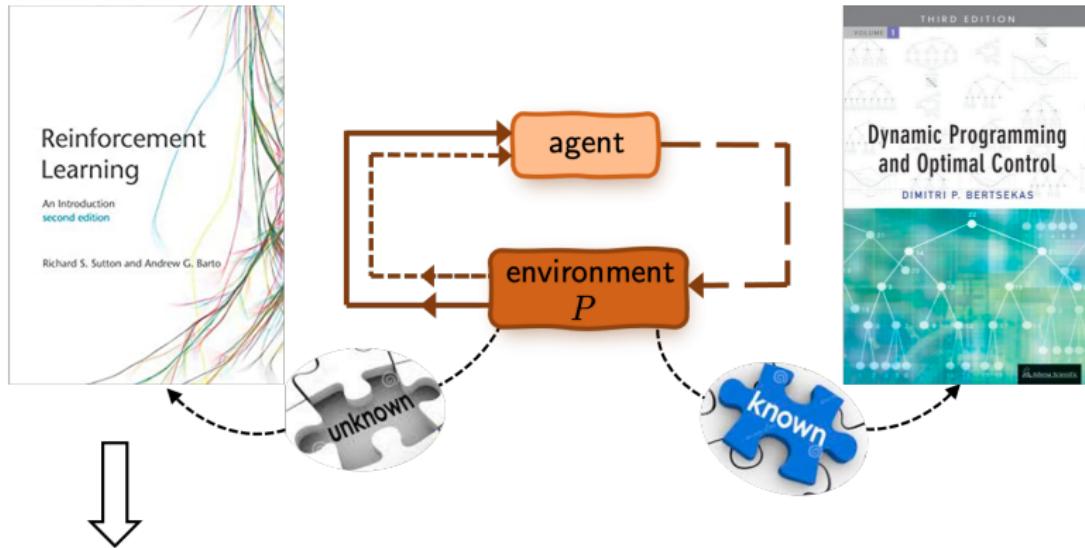
Yuantao Gu  
Tsinghua EE

# When the model is known . . .



**Planning:** computing the optimal policy  $\pi^*$  given MDP specification

# When the model is unknown ...

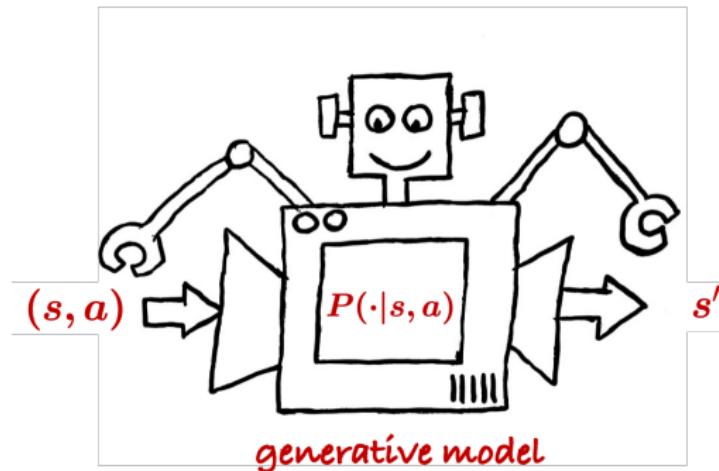


Need to learn optimal policy from samples w/o model specification

# This work: RL with a generative model / simulator

---

— Kearns, Singh '99



For each state-action pair  $(s, a)$ , collect  $N$  samples  $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$

**Question:** how many samples are sufficient to  
learn an  $\varepsilon$ -optimal policy ?

**Question:** how many samples are sufficient to learn an  $\varepsilon$ -optimal policy ?

$$\underbrace{\forall s: V^{\widehat{\pi}}(s) \geq V^*(s) - \varepsilon}$$

## An incomplete list of prior art

---

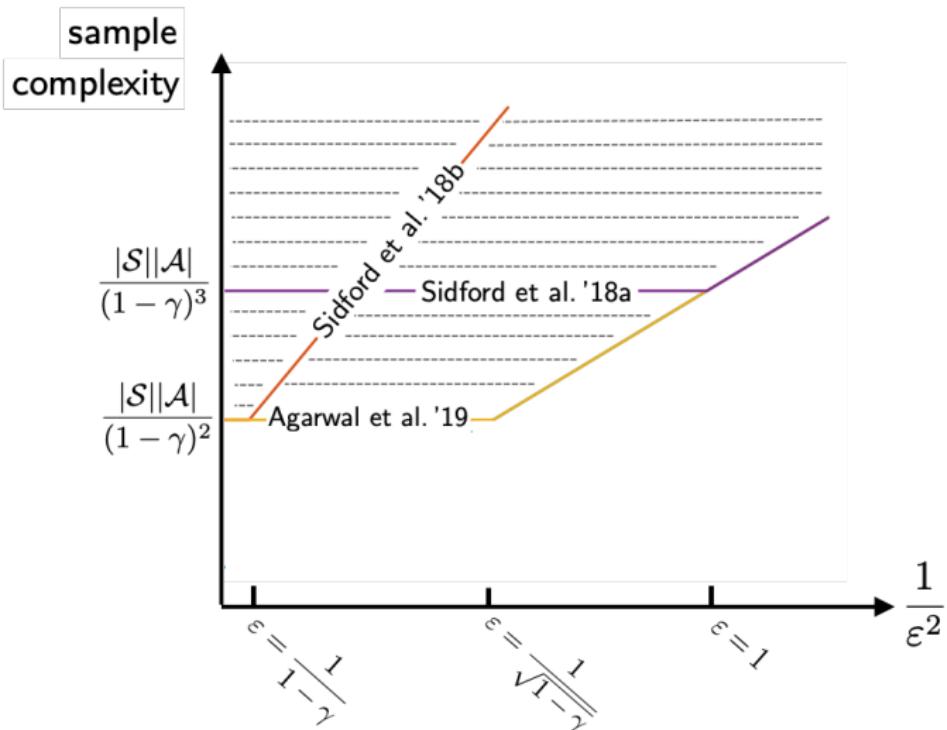
- Kearns & Singh '99
- Kakade '03
- Kearns, Mansour & Ng '02
- Azar, Munos & Kappen '12
- Azar, Munos, Ghavamzadeh & Kappen '13
- Sidford, Wang, Wu, Yang & Ye '18
- Sidford, Wang, Wu & Ye '18
- Wang '17
- Agarwal, Kakade & Yang '19
- Wainwright '19a
- Wainwright '19b
- Pananjady & Wainwright '20
- Yang & Wang '19
- Khamaru, Pananjady, Ruan, Wainwright & Jordan '20
- Mou, Li, Wainwright, Bartlett & Jordan '20
- ...

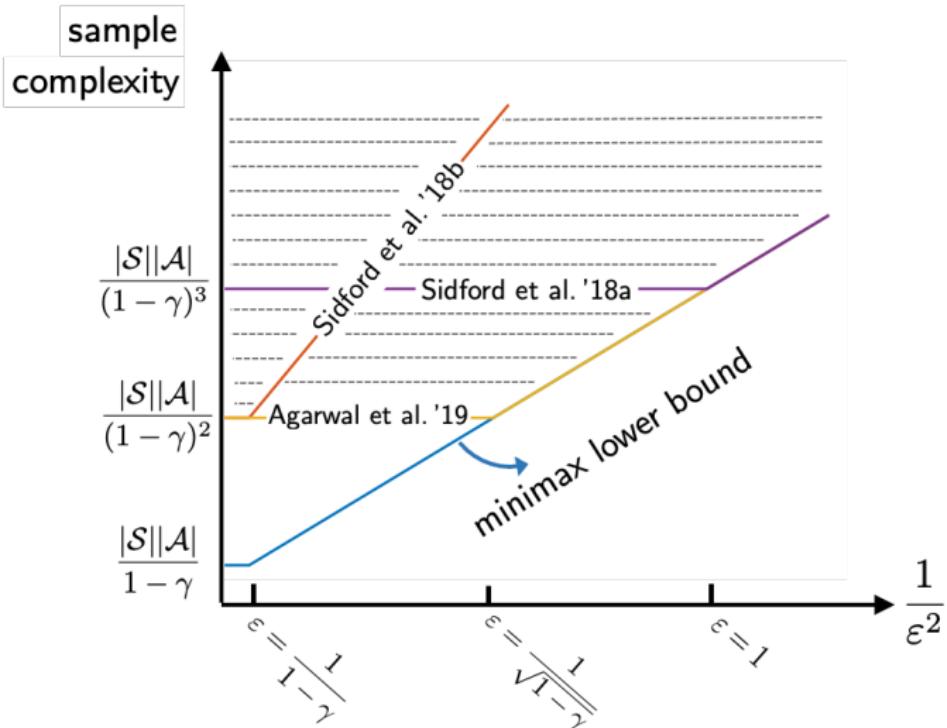
## An even shorter list of prior art

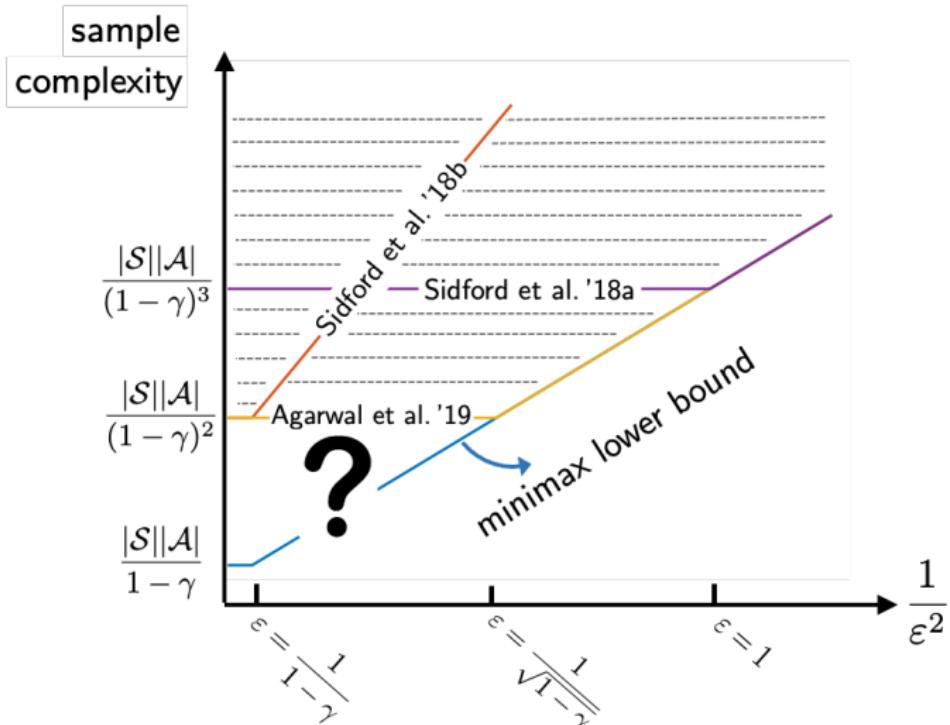
---

algorithm	sample size range	sample complexity	$\varepsilon$ -range
empirical QVI Azar et al. '13	$\left[ \frac{ \mathcal{S} ^2  \mathcal{A} }{(1-\gamma)^2}, \infty \right)$	$\frac{ \mathcal{S}  \mathcal{A} }{(1-\gamma)^3 \varepsilon^2}$	$(0, \frac{1}{\sqrt{(1-\gamma) \mathcal{S} }}]$
sublinear randomized VI Sidford et al. '18a	$\left[ \frac{ \mathcal{S}  \mathcal{A} }{(1-\gamma)^2}, \infty \right)$	$\frac{ \mathcal{S}  \mathcal{A} }{(1-\gamma)^4 \varepsilon^2}$	$(0, \frac{1}{1-\gamma}]$
variance-reduced QVI Sidford et al. '18b	$\left[ \frac{ \mathcal{S}  \mathcal{A} }{(1-\gamma)^3}, \infty \right)$	$\frac{ \mathcal{S}  \mathcal{A} }{(1-\gamma)^3 \varepsilon^2}$	$(0, 1]$
<b>empirical MDP + planning</b> Agarwal et al. '19	$\left[ \frac{ \mathcal{S}  \mathcal{A} }{(1-\gamma)^2}, \infty \right)$	$\frac{ \mathcal{S}  \mathcal{A} }{(1-\gamma)^3 \varepsilon^2}$	$(0, \frac{1}{\sqrt{1-\gamma}}]$

— see also Wainwright '19 (for estimating optimal values)





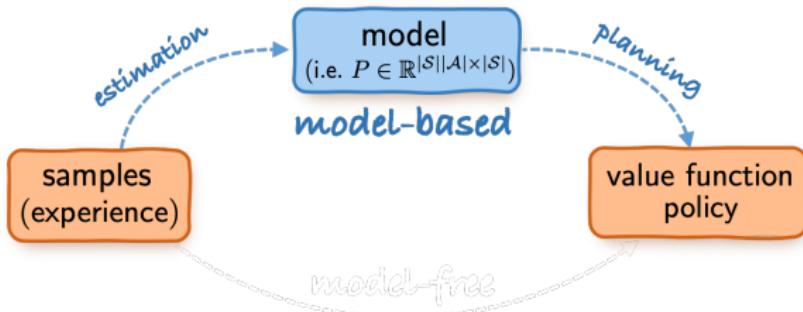


All prior theory requires sample size  $> \underbrace{\frac{|S||\mathcal{A}|}{(1 - \gamma)^2}}_{\text{sample size barrier}}$

*Is it possible to close the gap?*

# Two approaches

---

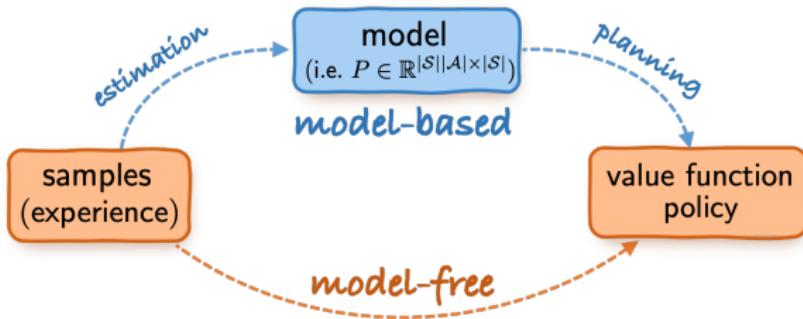


## Model-based approach (“plug-in”)

1. build an empirical estimate  $\hat{P}$  for  $P$
2. planning based on empirical  $\hat{P}$

# Two approaches

---



## Model-based approach (“plug-in”)

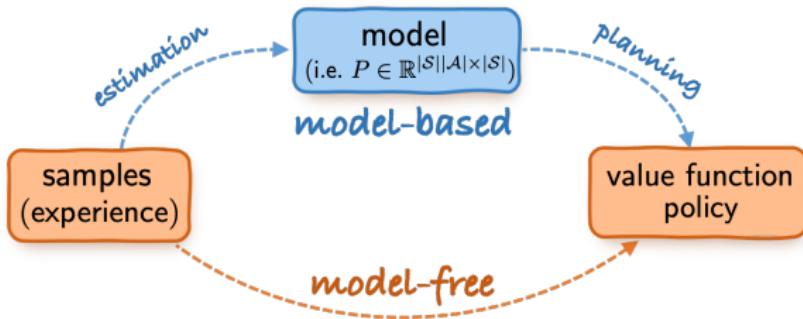
1. build an empirical estimate  $\hat{P}$  for  $P$
2. planning based on empirical  $\hat{P}$

## Model-free approach

— learning w/o constructing model explicitly

# Two approaches

---



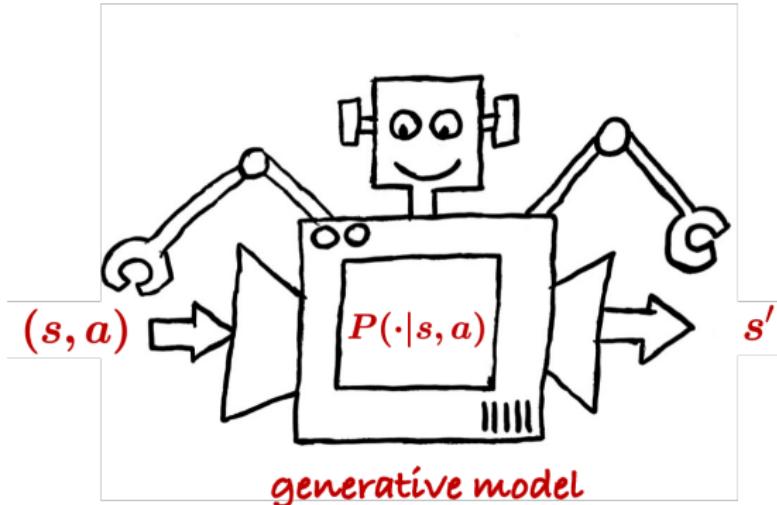
## Model-based approach (“plug-in”)

1. build empirical estimate  $\hat{P}$  for  $P$
2. planning based on empirical  $\hat{P}$

## Model-free approach

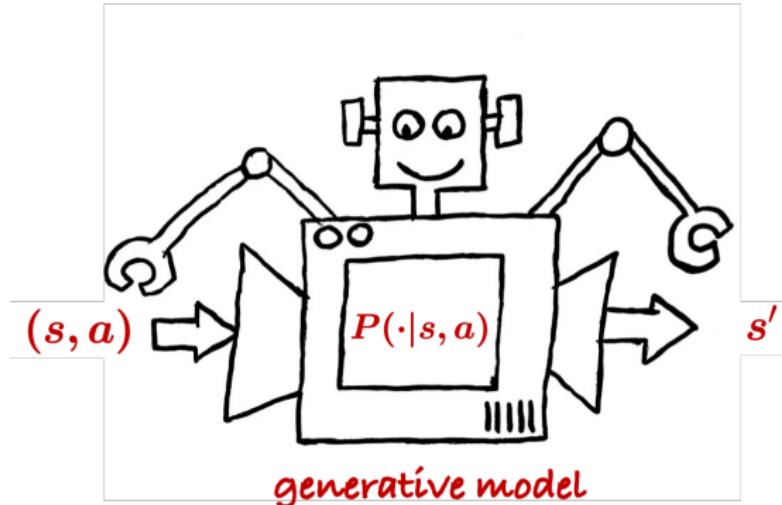
— learning w/o constructing model explicitly

# Model estimation



**Sampling:** for each  $(s, a)$ , collect  $N$  ind. samples  $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$

# Model estimation



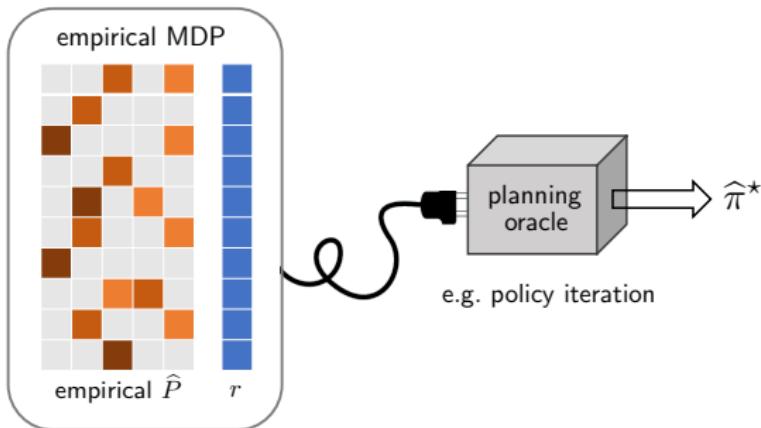
**Sampling:** for each  $(s, a)$ , collect  $N$  ind. samples  $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$

**Empirical estimates:** estimate  $\hat{P}(s'|s, a)$  by  $\underbrace{\frac{1}{N} \sum_{i=1}^N \mathbb{1}\{s'_{(i)} = s'\}}_{\text{empirical frequency}}$

# Model-based (plug-in) estimator

---

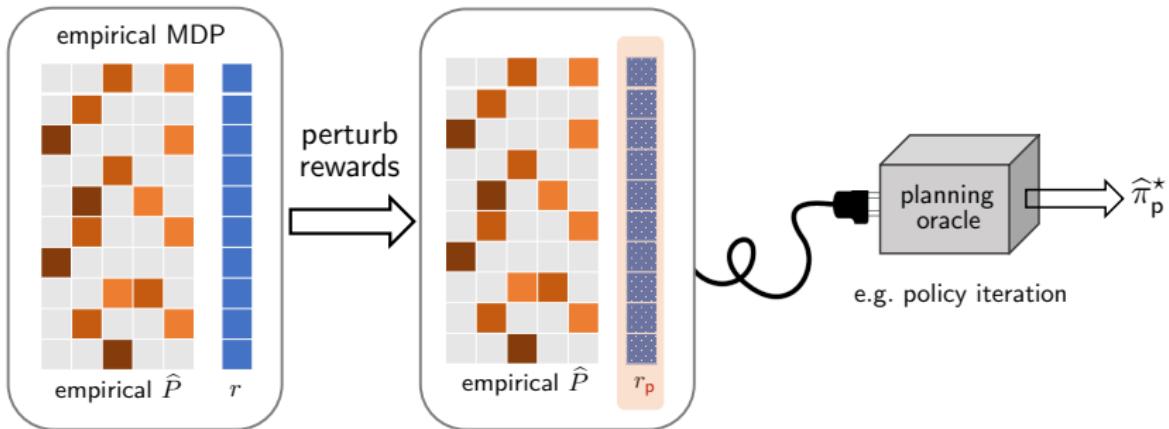
— Azar et al. '13, Agarwal et al. '19, Pananjady et al. '20



Planning based on the *empirical* MDP with *slightly perturbed rewards*

# Our method: plug-in estimator + perturbation

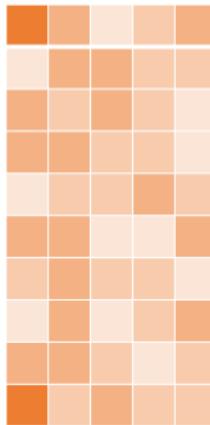
— Li, Wei, Chi, Gu, Chen '20



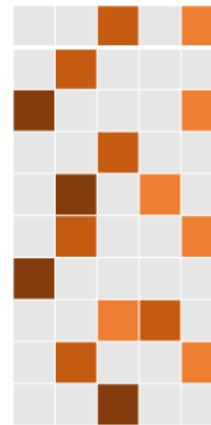
Run planning algorithms based on the *empirical* MDP

# Challenges in the sample-starved regime

---



truth:  
 $P \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$

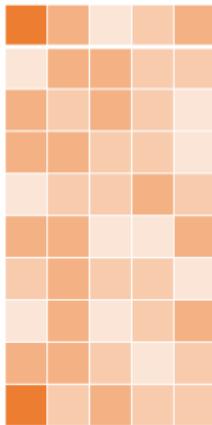


empirical estimate:  
 $\hat{P}$

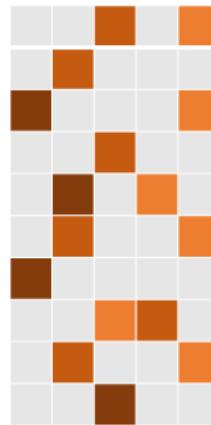
- Can't recover  $P$  faithfully if sample size  $\ll |\mathcal{S}|^2|\mathcal{A}|$ !

# Challenges in the sample-starved regime

---



truth:  
 $P \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$



empirical estimate:  
 $\hat{P}$

- Can't recover  $P$  faithfully if sample size  $\ll |\mathcal{S}|^2|\mathcal{A}|$ !
- Can we trust our policy estimate when reliable model estimation is infeasible?

# Main result

---

## Theorem 1 (Li, Wei, Chi, Gu, Chen '20)

For any  $0 < \varepsilon \leq \frac{1}{1-\gamma}$ , the optimal policy  $\widehat{\pi}_p^*$  of the perturbed empirical MDP achieves

$$\|V^{\widehat{\pi}_p^*} - V^*\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$$

# Main result

## Theorem 1 (Li, Wei, Chi, Gu, Chen '20)

For any  $0 < \varepsilon \leq \frac{1}{1-\gamma}$ , the optimal policy  $\hat{\pi}_p^*$  of the perturbed empirical MDP achieves

$$\|V^{\hat{\pi}_p^*} - V^*\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$$

- $\hat{\pi}_p^*$ : obtained by empirical QVI or PI within  $\tilde{O}\left(\frac{1}{1-\gamma}\right)$  iterations

# Main result

## Theorem 1 (Li, Wei, Chi, Gu, Chen '20)

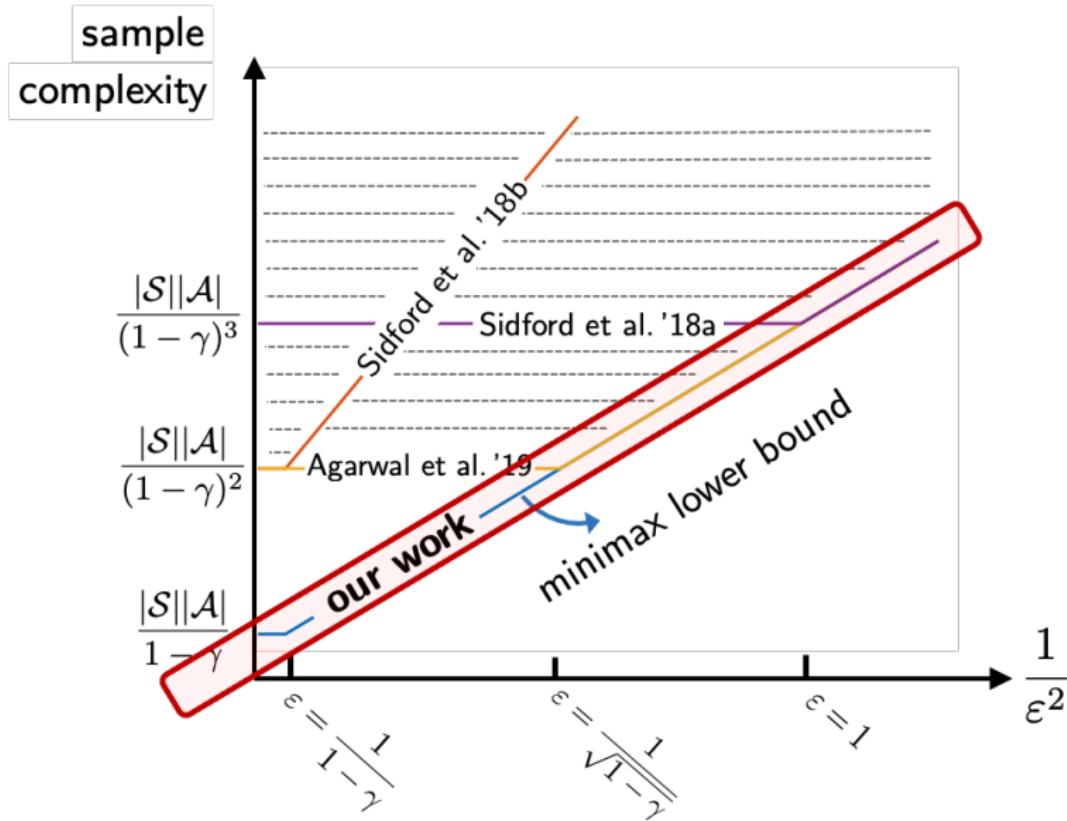
For any  $0 < \varepsilon \leq \frac{1}{1-\gamma}$ , the optimal policy  $\widehat{\pi}_p^*$  of the perturbed empirical MDP achieves

$$\|V^{\widehat{\pi}_p^*} - V^*\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$$

- $\widehat{\pi}_p^*$ : obtained by empirical QVI or PI within  $\tilde{O}(\frac{1}{1-\gamma})$  iterations
- **Minimax lower bound:**  $\tilde{\Omega}(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2})$  (Azar et al. '13)



*Analysis*

# Notation and Bellman equation

---

- $V^\pi$ : true value function under policy  $\pi$ 
  - Bellman equation:  $V^\pi = (I - P_\pi)^{-1}r$

# Notation and Bellman equation

---

- $V^\pi$ : true value function under policy  $\pi$ 
  - Bellman equation:  $V^\pi = (I - P_\pi)^{-1}r$
- $\hat{V}^\pi$ : estimate of value function under policy  $\pi$ 
  - Bellman equation:  $\hat{V}^\pi = (I - \hat{P}_\pi)^{-1}r$

# Notation and Bellman equation

---

- $V^\pi$ : true value function under policy  $\pi$ 
  - Bellman equation:  $V^\pi = (I - P_\pi)^{-1}r$
- $\hat{V}^\pi$ : estimate of value function under policy  $\pi$ 
  - Bellman equation:  $\hat{V}^\pi = (I - \hat{P}_\pi)^{-1}r$
- $\pi^*$ : optimal policy w.r.t. true value function
- $\hat{\pi}^*$ : optimal policy w.r.t. empirical value function

# Notation and Bellman equation

---

- $V^\pi$ : true value function under policy  $\pi$ 
  - Bellman equation:  $V^\pi = (I - P_\pi)^{-1}r$
- $\hat{V}^\pi$ : estimate of value function under policy  $\pi$ 
  - Bellman equation:  $\hat{V}^\pi = (I - \hat{P}_\pi)^{-1}r$
- $\pi^*$ : optimal policy w.r.t. true value function
- $\hat{\pi}^*$ : optimal policy w.r.t. empirical value function
- $V^* := V^{\pi^*}$ : optimal values under true models
- $\hat{V}^* := \hat{V}^{\hat{\pi}^*}$ : optimal values under empirical models

# Proof ideas

---

Elementary decomposition:

$$V^* - V^{\widehat{\pi}^*} = (V^* - \widehat{V}^{\pi^*}) + (\widehat{V}^{\pi^*} - \widehat{V}^{\widehat{\pi}^*}) + (\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*})$$

# Proof ideas

---

Elementary decomposition:

$$\begin{aligned} V^* - V^{\hat{\pi}^*} &= (V^* - \hat{V}^{\pi^*}) + (\hat{V}^{\pi^*} - \hat{V}^{\hat{\pi}^*}) + (\hat{V}^{\hat{\pi}^*} - V^{\hat{\pi}^*}) \\ &\leq (V^{\pi^*} - \hat{V}^{\pi^*}) + \textcolor{red}{0} + (\hat{V}^{\hat{\pi}^*} - V^{\hat{\pi}^*}) \end{aligned}$$

- **Step 1:** control  $V^\pi - \hat{V}^\pi$  for a fixed  $\pi$   
**(Bernstein inequality + high-order decomposition)**

# Proof ideas

---

Elementary decomposition:

$$\begin{aligned} V^* - V^{\hat{\pi}^*} &= (V^* - \hat{V}^{\pi^*}) + (\hat{V}^{\pi^*} - \hat{V}^{\hat{\pi}^*}) + (\hat{V}^{\hat{\pi}^*} - V^{\hat{\pi}^*}) \\ &\leq (V^{\pi^*} - \hat{V}^{\pi^*}) + \mathbf{0} + (\hat{V}^{\hat{\pi}^*} - V^{\hat{\pi}^*}) \end{aligned}$$

- **Step 1:** control  $V^\pi - \hat{V}^\pi$  for a fixed  $\pi$   
**(Bernstein inequality + high-order decomposition)**
- **Step 2:** extend it to control  $\hat{V}^{\hat{\pi}^*} - V^{\hat{\pi}^*}$  ( $\hat{\pi}^*$  depends on samples)  
**(decouple statistical dependency)**

# Step 1: improved theory for policy evaluation

---

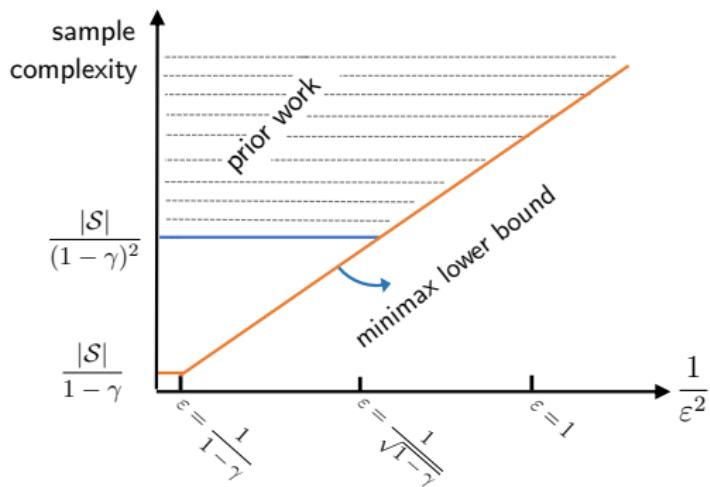
## Model-based policy evaluation:

- given a fixed policy  $\pi$ , estimate  $V^\pi$  via the plug-in estimate  $\hat{V}^\pi$

# Step 1: improved theory for policy evaluation

## Model-based policy evaluation:

— given a fixed policy  $\pi$ , estimate  $V^\pi$  via the plug-in estimate  $\hat{V}^\pi$



- A sample size barrier  $\frac{|\mathcal{S}|}{(1-\gamma)^2}$  already appeared in prior work  
(Agarwal et al. '19, Pananjady & Wainwright '19, Khamaru et al. '20)

# Step 1: improved theory for policy evaluation

## Model-based policy evaluation:

- given a fixed policy  $\pi$ , estimate  $V^\pi$  via the plug-in estimate  $\hat{V}^\pi$

### Theorem 2 (Li, Wei, Chi, Gu, Chen'20)

Fix any policy  $\pi$ . For  $0 < \varepsilon \leq \frac{1}{1-\gamma}$ , the plug-in estimator  $\hat{V}^\pi$  obeys

$$\|\hat{V}^\pi - V^\pi\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}|}{(1-\gamma)^3 \varepsilon^2}\right)$$

# Step 1: improved theory for policy evaluation

## Model-based policy evaluation:

— given a fixed policy  $\pi$ , estimate  $V^\pi$  via the plug-in estimate  $\hat{V}^\pi$

### Theorem 2 (Li, Wei, Chi, Gu, Chen'20)

Fix any policy  $\pi$ . For  $0 < \varepsilon \leq \frac{1}{1-\gamma}$ , the plug-in estimator  $\hat{V}^\pi$  obeys

$$\|\hat{V}^\pi - V^\pi\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}|}{(1-\gamma)^3\varepsilon^2}\right)$$

- Minimax optimal for all  $\varepsilon$  (Azar et al. '13, Pananjady & Wainwright '19)

## Key idea 1: a peeling argument

---

**Agarwal, Kakade, Yang 19:** first-order expansion

$$\hat{V}^\pi - V^\pi = \gamma(I - \gamma P_\pi)^{-1}(\hat{P}_\pi - P_\pi)\hat{V}^\pi \quad (\star)$$

**Ours:** higher-order expansion  $\longrightarrow$  tighter control

$$\hat{V}^\pi - V^\pi = \gamma(I - \gamma P_\pi)^{-1}(\hat{P}_\pi - P_\pi)\textcolor{red}{V}^\pi +$$

# Key idea 1: a peeling argument

---

**Agarwal, Kakade, Yang 19:** first-order expansion

$$\hat{V}^\pi - V^\pi = \gamma(I - \gamma P_\pi)^{-1}(\hat{P}_\pi - P_\pi)\hat{V}^\pi \quad (\star)$$

**Ours:** higher-order expansion  $\longrightarrow$  tighter control

$$\begin{aligned}\hat{V}^\pi - V^\pi &= \gamma(I - \gamma P_\pi)^{-1}(\hat{P}_\pi - P_\pi)\textcolor{red}{V}^\pi + \\ &\quad + \gamma(I - \gamma P_\pi)^{-1}(\hat{P}_\pi - P_\pi) \left( \hat{V}^\pi - V^\pi \right)\end{aligned}$$

# Key idea 1: a peeling argument

---

**Agarwal, Kakade, Yang 19:** first-order expansion

$$\widehat{V}^\pi - V^\pi = \gamma(I - \gamma P_\pi)^{-1}(\widehat{P}_\pi - P_\pi)\widehat{V}^\pi \quad (\star)$$

**Ours:** higher-order expansion  $\longrightarrow$  tighter control

$$\begin{aligned}\widehat{V}^\pi - V^\pi &= \gamma(I - \gamma P_\pi)^{-1}(\widehat{P}_\pi - P_\pi)\textcolor{red}{V}^\pi + \\ &\quad + \gamma^2 \left( (I - \gamma P_\pi)^{-1}(\widehat{P}_\pi - P_\pi) \right)^2 \textcolor{red}{V}^\pi \\ &\quad + \gamma^3 \left( (I - \gamma P_\pi)^{-1}(\widehat{P}_\pi - P_\pi) \right)^3 \textcolor{red}{V}^\pi \\ &\quad + \dots\end{aligned}$$

## Step 2: controlling $\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}$

---

A natural idea: apply our policy evaluation theory + union bound

## Step 2: controlling $\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}$

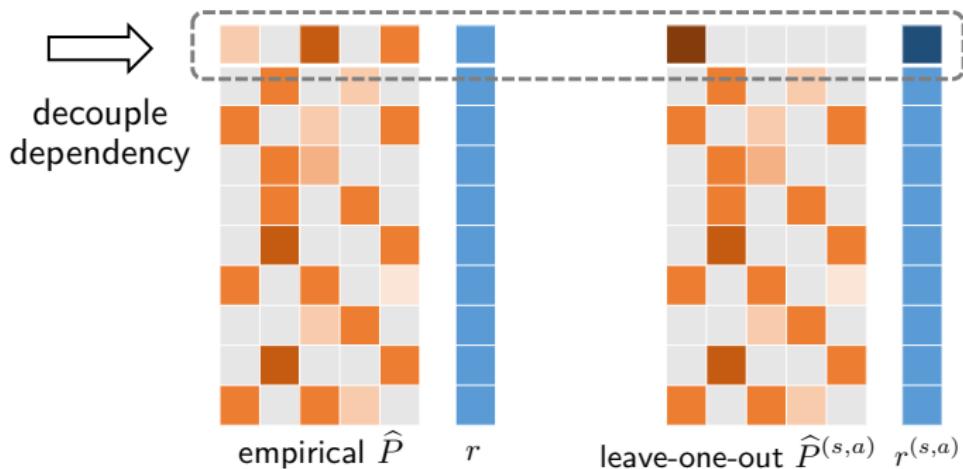
---

A natural idea: apply our policy evaluation theory + union bound

- highly suboptimal! (there are exponentially many policies)

## Key idea 2: leave-one-out analysis

Decouple dependency by introducing auxiliary state-action absorbing MDPs by dropping randomness for each  $(s, a)$



— inspired by Agarwal et al. '19 but quite different ...

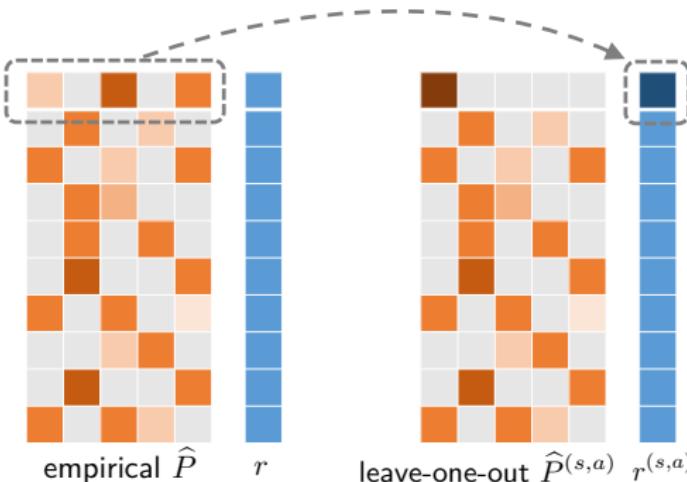
## Key idea 2: leave-one-out analysis

---

- Stein '72
- El Karoui, Bean, Bickel, Lim, Yu '13
- El Karoui '15
- Javanmard, Montanari '15
- Zhong, Boumal '17
- Lei, Bickel, El Karoui '17
- Sur, Chen, Candès '17
- Abbe, Fan, Wang, Zhong '17
- Chen, Fan, Ma, Wang '17
- Ma, Wang, Chi, Chen '17
- Chen, Chi, Fan, Ma '18
- Ding, Chen '18
- Dong, Shi '18
- Chen, Chi, Fan, Ma, Yan '19
- Chen, Fan, Ma, Yan '19
- Cai, Li, Poor, Chen '19
- Agarwal, Kakade, Yang '19
- Pananjady, Wainwright '19
- Ling '20

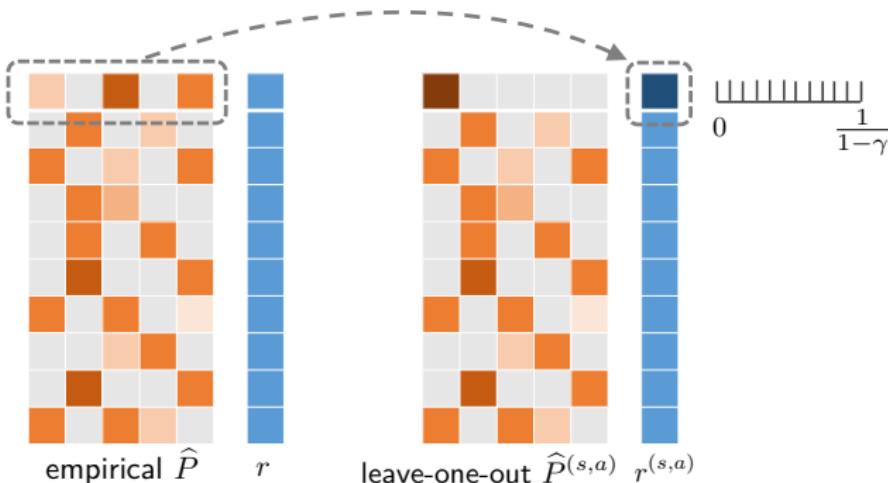
## Key idea 2: leave-one-out analysis

---



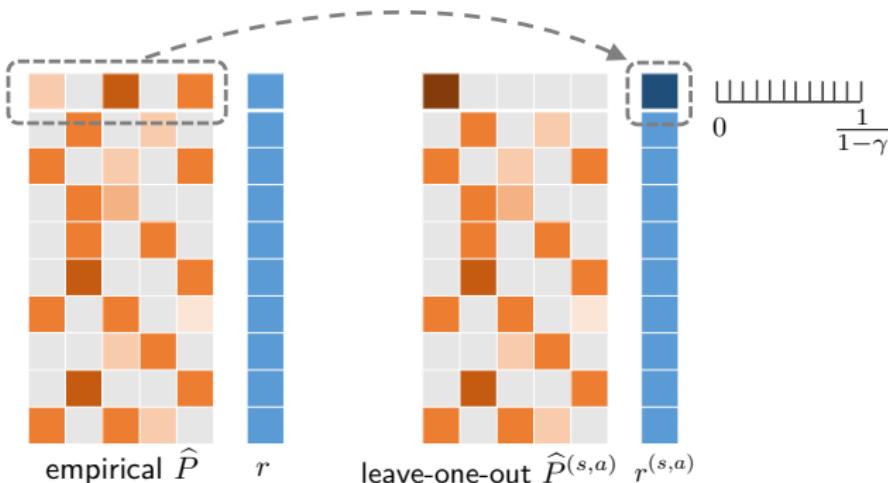
1. embed all randomness from  $\hat{P}_{s,a}$  into a single scalar (i.e.  $r_{s,a}^{(s,a)}$ )

## Key idea 2: leave-one-out analysis



1. embed all randomness from  $\hat{P}_{s,a}$  into a single scalar (i.e.  $r_{s,a}^{(s,a)}$ )
2. build an  $\epsilon$ -net for this scalar

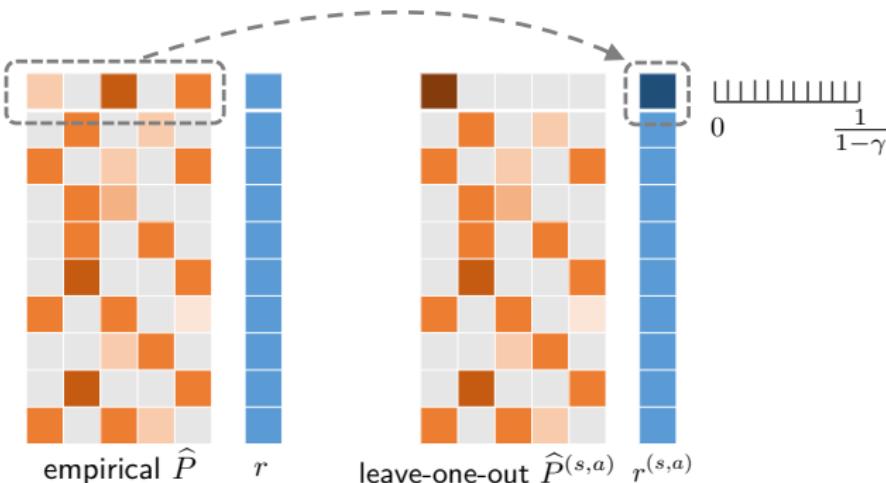
## Key idea 2: leave-one-out analysis



1. embed all randomness from  $\hat{P}_{s,a}$  into a single scalar (i.e.  $r_{s,a}^{(s,a)}$ )
2. build an  $\epsilon$ -net for this scalar
3.  $\hat{\pi}^*$  can be determined by this  $\epsilon$ -net under separation condition

$$\forall s \in \mathcal{S}, \quad \hat{Q}^*(s, \hat{\pi}^*(s)) - \max_{a: a \neq \hat{\pi}^*(s)} \hat{Q}^*(s, a) > 0$$

## Key idea 2: leave-one-out analysis



### Our decoupling argument vs. Agarwal, Kakade, Yang '19

- Agarwal et al. '19: dependency btw value  $\hat{V}$  & samples
- Ours: dependency btw policy  $\hat{\pi}$  & samples

## Key idea 3: tie-breaking via perturbation

---

- How to ensure separation between the optimal policy and others?

$$\forall s \in \mathcal{S}, \quad \hat{Q}^*(s, \hat{\pi}^*(s)) - \max_{a: a \neq \hat{\pi}^*(s)} \hat{Q}^*(s, a) > 0$$

## Key idea 3: tie-breaking via perturbation

---

- How to ensure separation between the optimal policy and others?

$$\forall s \in \mathcal{S}, \quad \hat{Q}^*(s, \hat{\pi}^*(s)) - \max_{a: a \neq \hat{\pi}^*(s)} \hat{Q}^*(s, a) > 0$$

- **Solution:** slightly perturb rewards  $r \implies \hat{\pi}_p^*$

- ensures  $\hat{\pi}_p^*$  can be differentiated from others
  - $V^{\hat{\pi}_p^*} \approx V^{\hat{\pi}^*}$



## Key idea 3: tie-breaking via perturbation

- How to ensure separation between the optimal policy and others?

$$\forall s \in \mathcal{S}, \quad \hat{Q}^*(s, \hat{\pi}^*(s)) - \max_{a: a \neq \hat{\pi}^*(s)} \hat{Q}^*(s, a) > \frac{(1-\gamma)\varepsilon}{|\mathcal{S}|^5 |\mathcal{A}|^5}$$

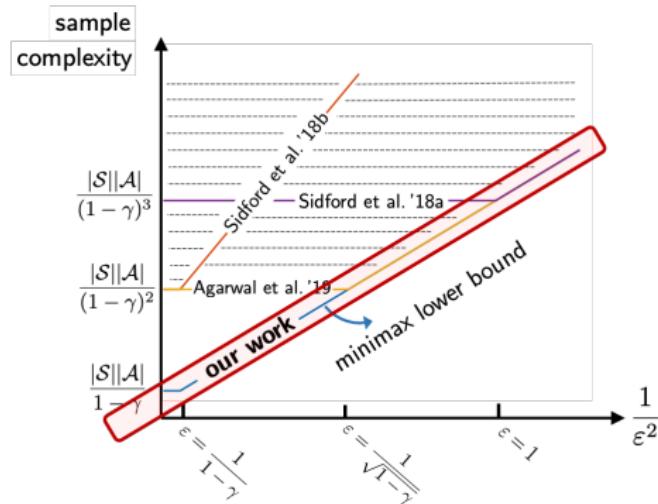
- **Solution:** slightly perturb rewards  $r$   $\implies \hat{\pi}_p^*$

- ensures  $\hat{\pi}_p^*$  can be differentiated from others
  - $V^{\hat{\pi}_p^*} \approx V^{\hat{\pi}^*}$



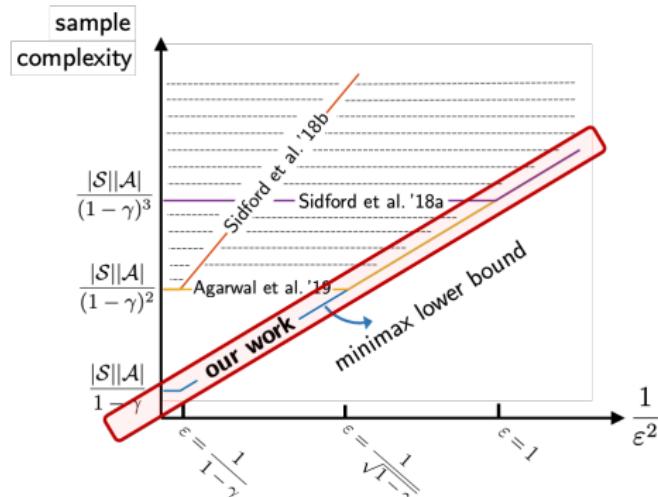
# Summary

Model-based RL is minimax optimal and does not suffer from a sample size barrier!



# Summary

Model-based RL is minimax optimal and does not suffer from a sample size barrier!



## future directions

- finite-horizon episodic MDPs
- Markov games

*Story 2: fast global convergence of entropy-regularized  
**natural policy gradient (NPG) methods***



Shicong Cen  
CMU ECE



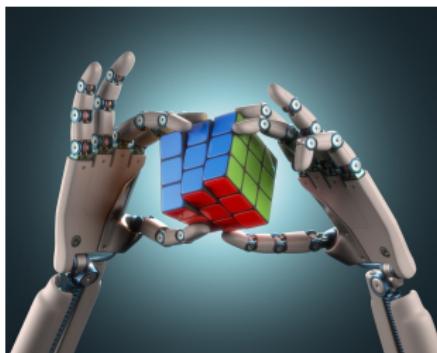
Chen Cheng  
Stanford Stats



Yuting Wei  
CMU Stats



Yuejie Chi  
CMU ECE



*Policy optimization: a major contributor to these successes*

# Policy gradient (PG) methods

---

Given initial state distribution  $s \sim \rho$ :

$$\text{maximize}_{\pi} \quad V^{\pi}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi}(s)]$$

# Policy gradient (PG) methods

---

Given initial state distribution  $s \sim \rho$ :

$$\text{maximize}_{\pi} \quad V^{\pi}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi}(s)]$$



**softmax parameterization:**

$$\pi_{\theta}(a|s) = \frac{\exp(\theta(s, a))}{\sum_a \exp(\theta(s, a))}$$

# Policy gradient (PG) methods

---

Given initial state distribution  $s \sim \rho$ :

$$\text{maximize}_{\pi} \quad V^{\pi}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi}(s)]$$



**softmax parameterization:**

$$\pi_{\theta}(a|s) = \frac{\exp(\theta(s, a))}{\sum_a \exp(\theta(s, a))}$$

$$\text{maximize}_{\theta} \quad V^{\pi_{\theta}}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi_{\theta}}(s)]$$

# Policy gradient (PG) methods

Given initial state distribution  $s \sim \rho$ :

$$\text{maximize}_{\pi} \quad V^{\pi}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi}(s)]$$



**softmax parameterization:**

$$\pi_{\theta}(a|s) = \frac{\exp(\theta(s, a))}{\sum_a \exp(\theta(s, a))}$$

$$\text{maximize}_{\theta} \quad V^{\pi_{\theta}}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi_{\theta}}(s)]$$

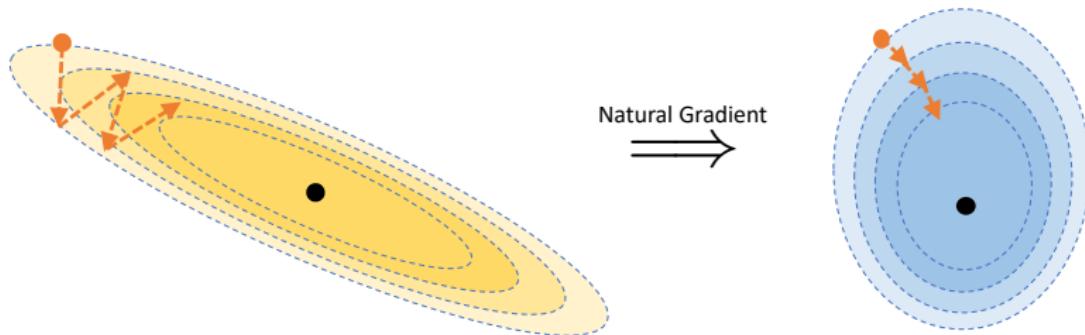
PG method (Sutton et al. '00)

$$\theta^{(t+1)} = \theta^{(t)} + \eta \nabla_{\theta} V^{\pi_{\theta}^{(t)}}(\rho), \quad t = 0, 1, \dots$$

- $\eta$ : learning rate

# Booster 1: natural policy gradient (NPG)

*precondition gradients to improve search directions ...*



NPG method (Kakade '02)

$$\theta^{(t+1)} = \theta^{(t)} + \eta (\mathcal{F}_\rho^\theta)^\dagger \nabla_\theta V^{\pi_\theta^{(t)}}(\rho), \quad t = 0, 1, \dots$$

- $\mathcal{F}_\rho^\theta := \mathbb{E} \left[ (\nabla_\theta \log \pi_\theta(a|s)) (\nabla_\theta \log \pi_\theta(a|s))^\top \right]$ : Fisher info matrix

## Booster 2: entropy regularization

---

*accelerate convergence by regularizing objective function*

$$\begin{aligned} V_\tau^\pi(s_0) &:= \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t (r_t - \color{blue}{\tau \log \pi(a_t|s_t)}) \mid s_0 \right] \\ &= V^\pi(s) + \frac{\color{blue}{\tau}}{1-\gamma} \mathbb{E}_{s \sim d_s^\pi} \underbrace{\left[ - \sum_a \pi(a|s) \log \pi(a|s) \mid s_0 \right]}_{\text{entropy}} \end{aligned}$$

- $\tau$ : regularization parameter
- $d_s^\pi$ : discounted state visitation distribution

## Booster 2: entropy regularization

accelerate convergence by regularizing objective function

$$\begin{aligned} V_\tau^\pi(s_0) &:= \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t (r_t - \tau \log \pi(a_t|s_t)) \mid s_0 \right] \\ &= V^\pi(s) + \frac{\tau}{1-\gamma} \mathbb{E}_{s \sim d_s^\pi} \left[ \underbrace{- \sum_a \pi(a|s) \log \pi(a|s)}_{\text{entropy}} \mid s_0 \right] \end{aligned}$$

- $\tau$ : regularization parameter
- $d_s^\pi$ : discounted state visitation distribution

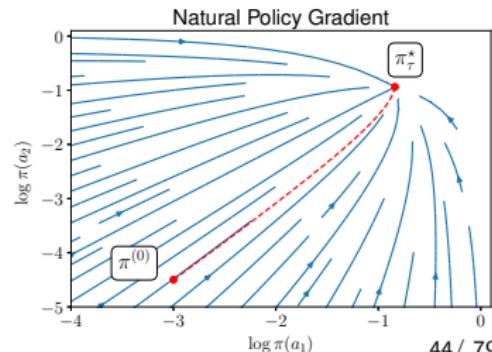
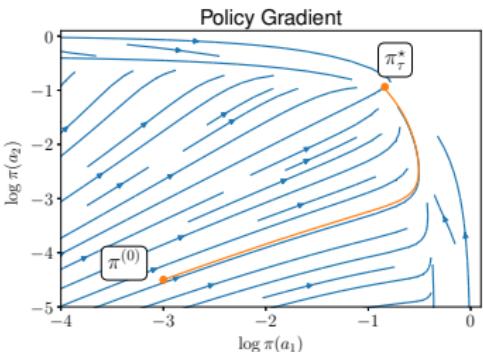
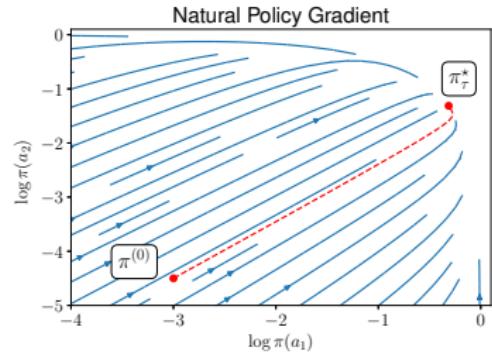
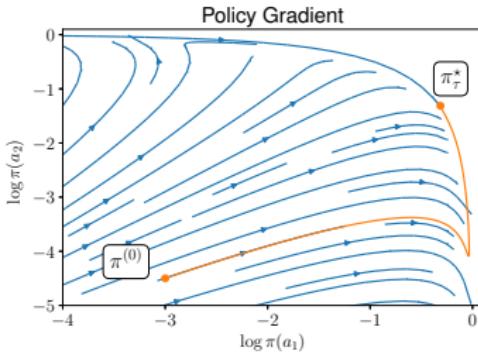
entropy-regularized value maximization

$$\text{maximize}_\theta \quad V_\tau^{\pi_\theta}(\rho) := \mathbb{E}_{s \sim \rho} [V_\tau^{\pi_\theta}(s)] \quad (\text{"soft" value function})$$

# Entropy-regularized natural gradient helps!

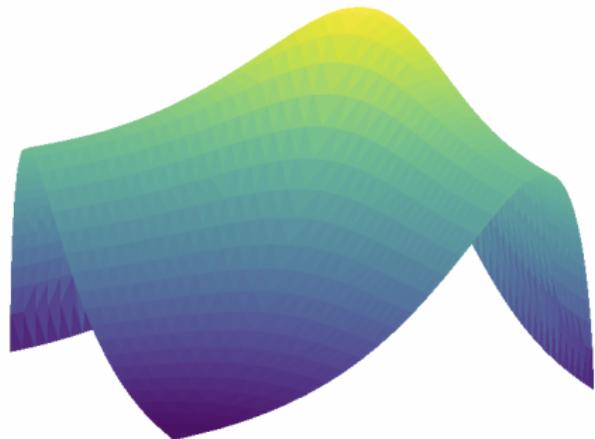
A toy bandit example: 3 arms with rewards 1, 0.9 and 0.1

increase regularization



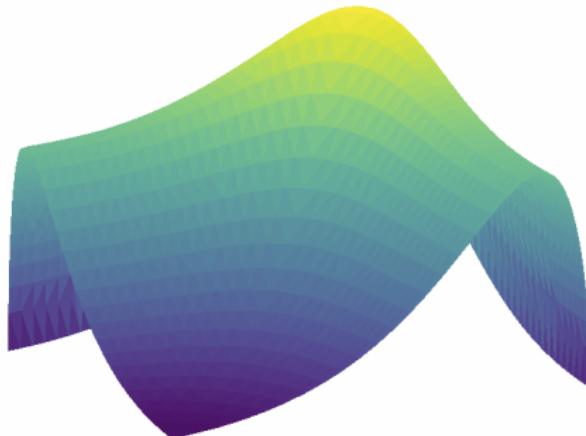
## Challenge: non-concavity

---



# Challenge: non-concavity

---

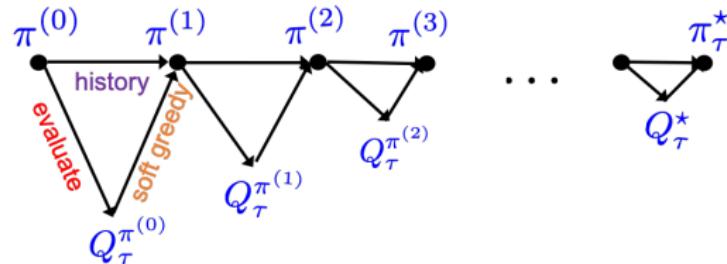


Recent advances

- PG for control ([Fazel et al., 2018; Bhandari and Russo, 2019](#))
- PG for tabular MDPs ([Agarwal et al. 19, Bhandari and Russo '19, Mei et al '20](#))
- unregularized NPG for tabular MDPs ([Agarwal et al. '19, Bhandari and Russo '20](#))
- ...

*This work: understanding entropy-regularized  
NPG methods in tabular settings*

# Entropy-regularized NPG in tabular settings



An alternative expression in policy space (tabular setting)

$$\pi^{(t+1)}(a|s) \propto \pi^{(t)}(a|s)^{1-\frac{\eta\tau}{1-\gamma}} \exp\left(\frac{\eta Q_\tau^{(t)}(s, a)}{1-\gamma}\right), \quad t = 0, 1, \dots$$

- $Q_\tau^{(t)}$ : soft Q-function of  $\pi^{(t)}$ ;  $0 < \eta \leq \frac{1-\gamma}{\tau}$ : learning rate

- invariant to the choice of initial state distribution  $\rho$

# Linear convergence with exact gradients

---

*optimal policy:*  $\pi_\tau^*$ ; *optimal “soft” Q function:*  $Q_\tau^* := Q_\tau^{\pi_\tau^*}$

**Exact oracle:** perfect gradient evaluation

## Theorem 3 (Cen, Cheng, Chen, Wei, Chi '20)

For any  $0 < \eta \leq (1 - \gamma)/\tau$ , entropy-regularized NPG achieves

$$\|Q_\tau^* - Q_\tau^{(t+1)}\|_\infty \leq C_1 \gamma (1 - \eta \tau)^t, \quad t = 0, 1, \dots$$

$$\bullet C_1 = \|Q_\tau^* - Q_\tau^{(0)}\|_\infty + 2\tau \left(1 - \frac{\eta\tau}{1-\gamma}\right) \|\log \pi_\tau^* - \log \pi^{(0)}\|_\infty$$

## Implications: iteration complexity

---

number of iterations needed to reach  $\|Q_\tau^* - Q_\tau^{(t+1)}\|_\infty \leq \varepsilon$  is at most

- **General learning rates** ( $0 < \eta < \frac{1-\gamma}{\tau}$ ):

$$\frac{1}{\eta\tau} \log \left( \frac{C_1 \gamma}{\varepsilon} \right)$$

- **Soft policy iteration** ( $\eta = \frac{1-\gamma}{\tau}$ ):

$$\frac{1}{1-\gamma} \log \left( \frac{\|Q_\tau^* - Q_\tau^{(0)}\|_\infty \gamma}{\varepsilon} \right)$$

## Implications: iteration complexity

---

number of iterations needed to reach  $\|Q_\tau^* - Q_\tau^{(t+1)}\|_\infty \leq \varepsilon$  is at most

- **General learning rates** ( $0 < \eta < \frac{1-\gamma}{\tau}$ ):

$$\frac{1}{\eta\tau} \log \left( \frac{C_1 \gamma}{\varepsilon} \right)$$

- **Soft policy iteration** ( $\eta = \frac{1-\gamma}{\tau}$ ):

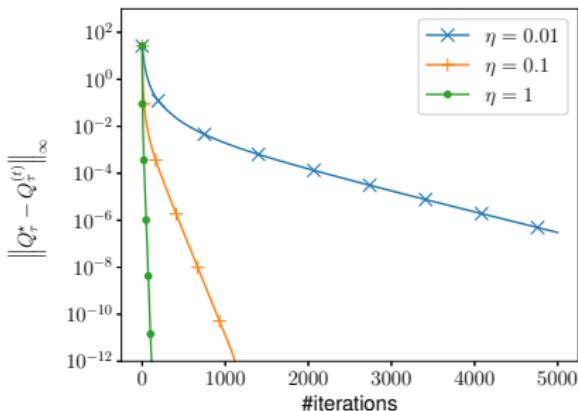
$$\frac{1}{1-\gamma} \log \left( \frac{\|Q_\tau^* - Q_\tau^{(0)}\|_\infty \gamma}{\varepsilon} \right)$$

Nearly dimension-free global linear convergence!

# Regularized NPG vs. unregularized NPG

regularized NPG

$\tau = 0.001$

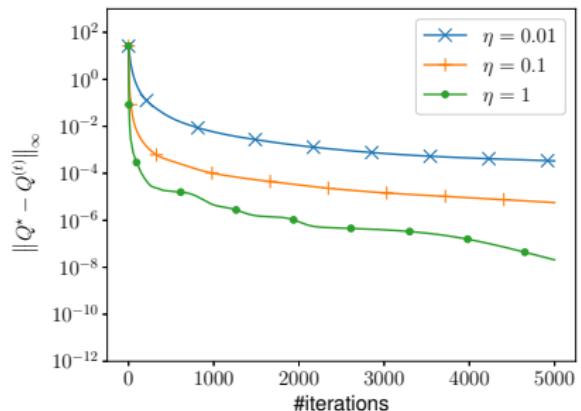


$$\text{linear rate: } \frac{1}{\eta\tau} \log\left(\frac{1}{\varepsilon}\right)$$

**Ours**

unregularized NPG

$\tau = 0$



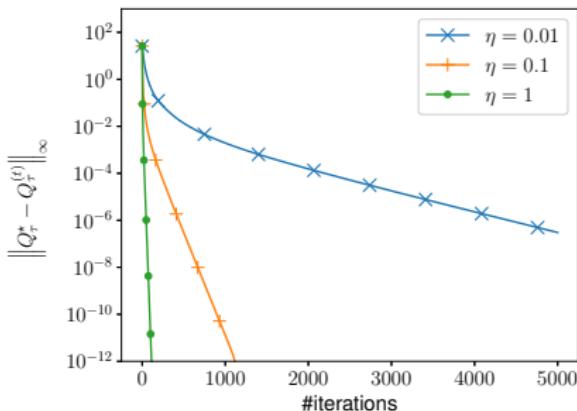
$$\text{sublinear rate: } \frac{1}{\min\{\eta, (1-\gamma)^2\}\varepsilon}$$

(Agarwal et al. '19)

# Regularized NPG vs. unregularized NPG

regularized NPG

$\tau = 0.001$

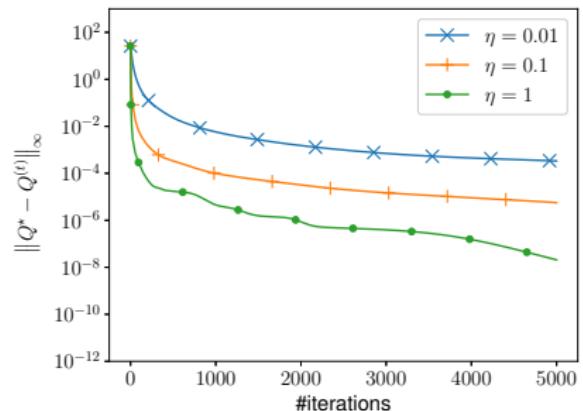


$$\text{linear rate: } \frac{1}{\eta\tau} \log\left(\frac{1}{\varepsilon}\right)$$

**Ours**

unregularized NPG

$\tau = 0$



$$\text{sublinear rate: } \frac{1}{\min\{\eta, (1-\gamma)^2\}\varepsilon}$$

(Agarwal et al. '19)

Entropy regularization enables faster convergence!

## Returning to the original MDP?

---

How to employ entropy-regularized NPG to find an  $\varepsilon$ -optimal policy for the original (unregularized) MDP?

- suffices to find an  $\frac{\varepsilon}{2}$ -optimal policy of regularized MDP  
w/ regularization parameter  $\tau = \frac{(1-\gamma)\varepsilon}{4 \log |\mathcal{A}|}$
- iteration complexity is the same as before (up to log factor)

# Entropy-regularized NPG with inexact gradients

---

**Inexact oracle:** inexact evaluation of  $Q_\tau^{(t)}$ , which returns  $\hat{Q}_\tau^{(t)}$  s.t.

$$\|\hat{Q}_\tau^{(t)} - Q_\tau^{(t)}\|_\infty \leq \delta,$$

e.g. using sample-based estimators

# Entropy-regularized NPG with inexact gradients

---

**Inexact oracle:** inexact evaluation of  $Q_\tau^{(t)}$ , which returns  $\hat{Q}_\tau^{(t)}$  s.t.

$$\|\hat{Q}_\tau^{(t)} - Q_\tau^{(t)}\|_\infty \leq \delta,$$

e.g. using sample-based estimators

**Inexact entropy-regularized NPG:**

$$\pi^{(t+1)}(a|s) \propto (\pi^{(t)}(a|s))^{1-\frac{\eta\tau}{1-\gamma}} \exp\left(\frac{\eta \hat{Q}_\tau^{(t)}(s, a)}{1-\gamma}\right)$$

# Entropy-regularized NPG with inexact gradients

---

**Inexact oracle:** inexact evaluation of  $Q_\tau^{(t)}$ , which returns  $\hat{Q}_\tau^{(t)}$  s.t.

$$\|\hat{Q}_\tau^{(t)} - Q_\tau^{(t)}\|_\infty \leq \delta,$$

e.g. using sample-based estimators

**Inexact entropy-regularized NPG:**

$$\pi^{(t+1)}(a|s) \propto (\pi^{(t)}(a|s))^{1-\frac{\eta\tau}{1-\gamma}} \exp\left(\frac{\eta \hat{Q}_\tau^{(t)}(s, a)}{1-\gamma}\right)$$

**Question:** stability vis-à-vis inexact gradient evaluation?

# Linear convergence with inexact gradients

$$\|\hat{Q}_\tau^{(t)} - Q_\tau^{(t)}\|_\infty \leq \delta$$

## Theorem 4 (Cen, Cheng, Chen, Wei, Chi '20)

For any stepsize  $0 < \eta \leq (1 - \gamma)/\tau$ , entropy-regularized NPG attains

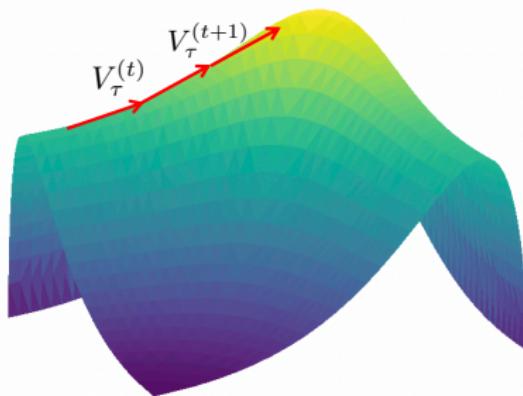
$$\|Q_\tau^\star - Q_\tau^{(t+1)}\|_\infty \leq \gamma(1 - \eta\tau)^t C_1 + C_2$$

- $C_1 = \|Q_\tau^\star - Q_\tau^{(0)}\|_\infty + 2\tau \left(1 - \frac{\eta\tau}{1 - \gamma}\right) \|\log \pi_\tau^\star - \log \pi^{(0)}\|_\infty$
- $C_2 = \frac{2\gamma \left(1 + \frac{\gamma}{\eta\tau}\right)}{1 - \gamma} \delta$ : error floor
- converges linearly at the same rate until an error floor is hit

*A little analysis when  $\eta = \frac{1-\gamma}{\tau}$*

# A key lemma: monotonic performance improvement

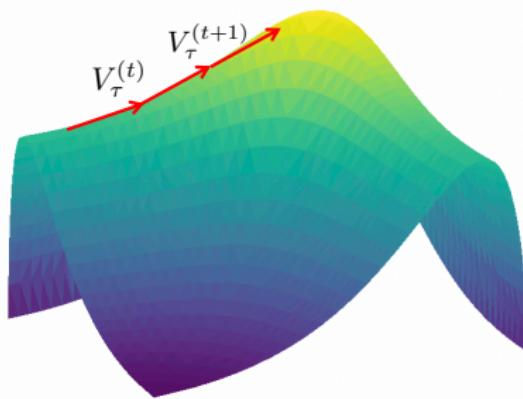
---



$$V_\tau^{(t+1)}(\rho) - V_\tau^{(t)}(\rho) = \mathbb{E}_{s \sim d_\rho^{(t+1)}} \left[ \underbrace{\left( \frac{1}{\eta} - \frac{\tau}{1-\gamma} \right) \text{KL}\left( \pi^{(t+1)}(\cdot|s) \parallel \pi^{(t)}(\cdot|s) \right)}_{\text{KL divergence}} + \underbrace{\frac{1}{\eta} \text{KL}\left( \pi^{(t)}(\cdot|s) \parallel \pi^{(t+1)}(\cdot|s) \right)}_{\text{KL divergence}} \right]$$

# A key lemma: monotonic performance improvement

---



$$\begin{aligned} V_\tau^{(t+1)}(\rho) - V_\tau^{(t)}(\rho) &= \mathbb{E}_{s \sim d_\rho^{(t+1)}} \left[ \underbrace{\left( \frac{1}{\eta} - \frac{\tau}{1-\gamma} \right) \text{KL}\left( \pi^{(t+1)}(\cdot|s) \parallel \pi^{(t)}(\cdot|s) \right)}_{\text{KL divergence}} \right. \\ &\quad \left. + \underbrace{\frac{1}{\eta} \text{KL}\left( \pi^{(t)}(\cdot|s) \parallel \pi^{(t+1)}(\cdot|s) \right)}_{\text{KL divergence}} \right] \\ &\geq 0 \quad (\text{if } 0 < \eta \leq \frac{1-\gamma}{\tau}) \end{aligned}$$

# “Soft” Bellman operator

---

$$\begin{aligned}\mathcal{T}_\tau(Q)(s, a) := & \underbrace{r(s, a)}_{\text{immediate reward}} \\ & + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \max_{\pi(\cdot|s')} \mathbb{E}_{a' \sim \pi(\cdot|s')} \left[ \underbrace{Q(s', a')}_{\text{next state's value}} - \underbrace{\tau \log \pi(a'|s')}_{\text{regularizer}} \right] \right]\end{aligned}$$

# “Soft” Bellman operator

---

$$\begin{aligned}\mathcal{T}_\tau(Q)(s, a) := & \underbrace{r(s, a)}_{\text{immediate reward}} \\ & + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \max_{\pi(\cdot|s')} \mathbb{E}_{a' \sim \pi(\cdot|s')} \left[ \underbrace{Q(s', a')}_{\text{next state's value}} - \underbrace{\tau \log \pi(a'|s')}_{\text{regularizer}} \right] \right]\end{aligned}$$

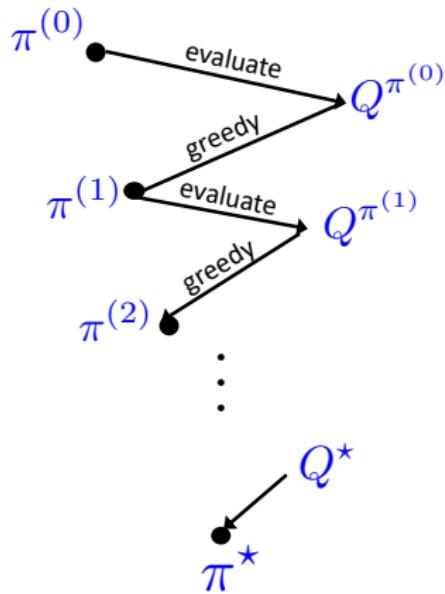
**Soft Bellman equation:**  $Q_\tau^*$  is the *unique* solution to

$$\mathcal{T}_\tau(Q) = Q$$

**$\gamma$ -contraction of soft Bellman operator:**

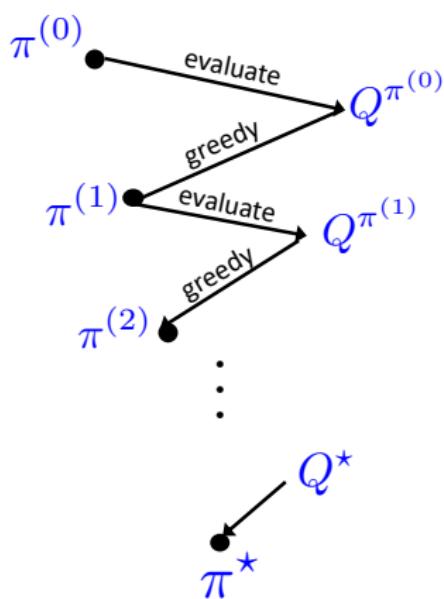
$$\|\mathcal{T}_\tau(Q_1) - \mathcal{T}_\tau(Q_2)\|_\infty \leq \gamma \|Q_1 - Q_2\|_\infty$$

## policy iteration



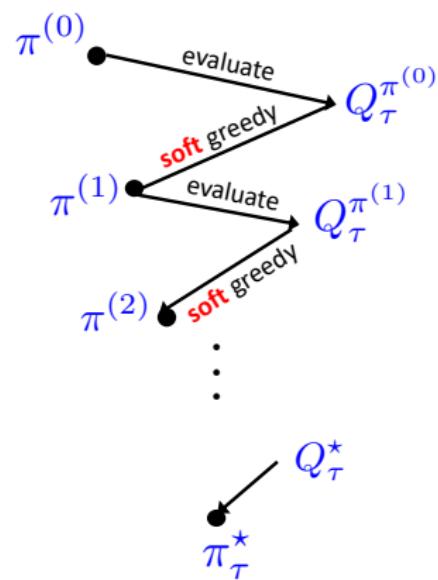
Bellman operator

policy iteration



Bellman operator

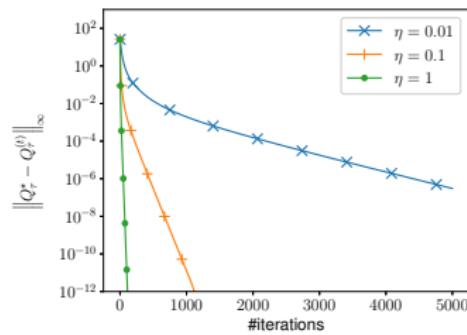
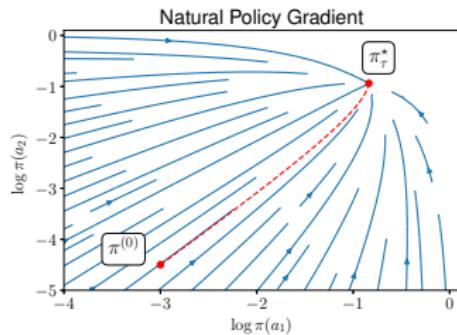
soft policy iteration ( $\eta = \frac{1-\gamma}{\tau}$ )



soft Bellman operator

# Summary

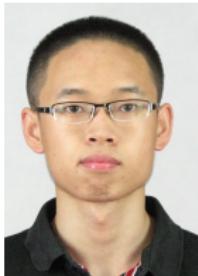
Global linear convergence of entropy-regularized NPG methods for tabular discounted MDPs



## future directions:

- function approximation
- sample complexities
- soft actor-critic algorithms

*Story 3: sample complexity of  
(asynchronous) Q-learning on Markovian samples*



Gen Li  
Tsinghua EE



Yuting Wei  
CMU Stats

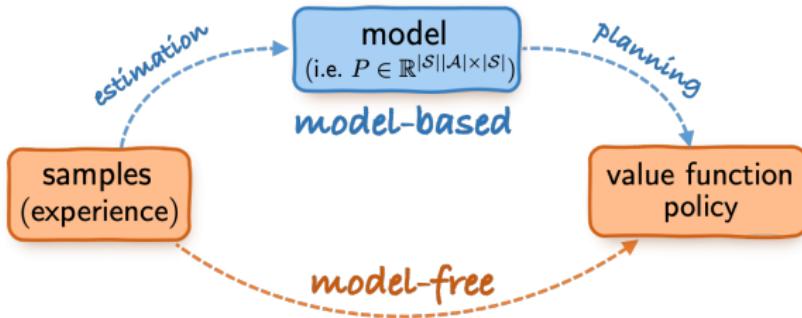


Yuejie Chi  
CMU ECE



Yuantao Gu  
Tsinghua EE

# Model-based vs. model-free RL

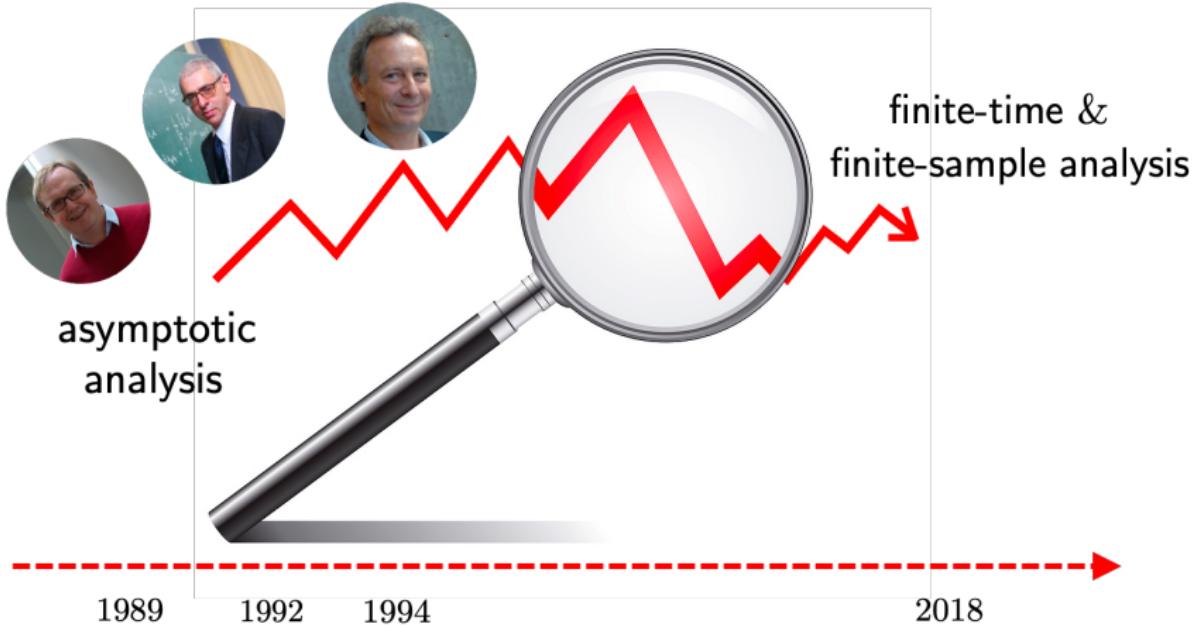


## Model-based approach (“plug-in”)

1. build an empirical estimate  $\hat{P}$  for  $P$
2. planning based on empirical  $\hat{P}$

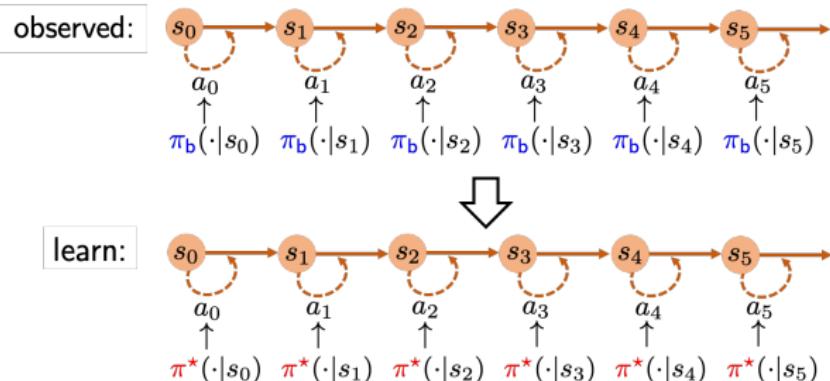
## Model-free approach

— learning w/o modeling & estimating environment explicitly



A classical example: **Q-learning** on Markovian samples

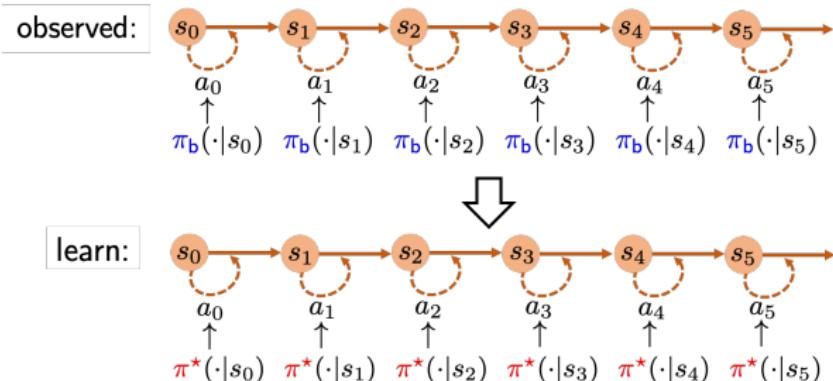
# Markovian samples and behavior policy



**Observed:**  $\underbrace{\{s_t, a_t, r_t\}_{t \geq 0}}_{\text{Markovian trajectory}}$  generated by behavior policy  $\pi_b$

**Goal:** learn optimal value  $V^*$  and  $Q^*$  based on sample trajectory

# Markovian samples and behavior policy



Key quantities of sample trajectory

- minimum state-action occupancy probability

$$\mu_{\min} := \min \underbrace{\mu_{\pi_b}(s, a)}_{\text{stationary distribution}}$$

- mixing time:  $t_{\text{mix}}$

# Q-learning: a classical model-free algorithm

---



Chris Watkins



Peter Dayan

Stochastic approximation for solving **Bellman equation**  $Q = \mathcal{T}(Q)$

 Robbins & Monro '51

## Aside: Bellman optimality principle

---

### Bellman operator

$$\mathcal{T}(Q)(s, a) := \underbrace{r(s, a)}_{\text{immediate reward}} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ \underbrace{\max_{a' \in \mathcal{A}} Q(s', a')}_{\text{next state's value}} \right]$$

- one-step look-ahead

## Aside: Bellman optimality principle

---

### Bellman operator

$$\mathcal{T}(Q)(s, a) := \underbrace{r(s, a)}_{\text{immediate reward}} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ \underbrace{\max_{a' \in \mathcal{A}} Q(s', a')}_{\text{next state's value}} \right]$$

- one-step look-ahead

**Bellman equation:**  $Q^*$  is *unique* solution to

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



*Richard Bellman*

# Q-learning: a classical model-free algorithm

---



Chris Watkins



Peter Dayan

Stochastic approximation for solving Bellman equation  $Q = \mathcal{T}(Q)$

$$\underbrace{Q_{t+1}(s_t, a_t) = (1 - \eta_t)Q_t(s_t, a_t) + \eta_t \mathcal{T}_t(Q_t)(s_t, a_t),}_{\text{only update } (s_t, a_t)\text{-th entry}} \quad t \geq 0$$

# Q-learning: a classical model-free algorithm

---



Chris Watkins



Peter Dayan

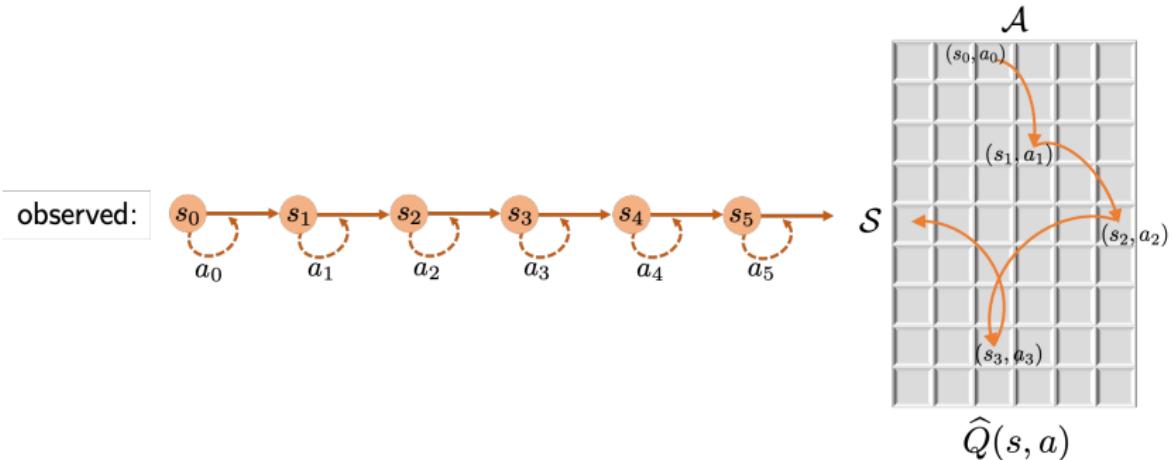
Stochastic approximation for solving Bellman equation  $Q = \mathcal{T}(Q)$

$$\underbrace{Q_{t+1}(s_t, a_t) = (1 - \eta_t)Q_t(s_t, a_t) + \eta_t \mathcal{T}_t(Q_t)(s_t, a_t),}_{\text{only update } (s_t, a_t)\text{-th entry}} \quad t \geq 0$$

$$\mathcal{T}_t(Q)(s_t, a_t) = r(s_t, a_t) + \gamma \max_{a'} Q(s_{t+1}, a')$$

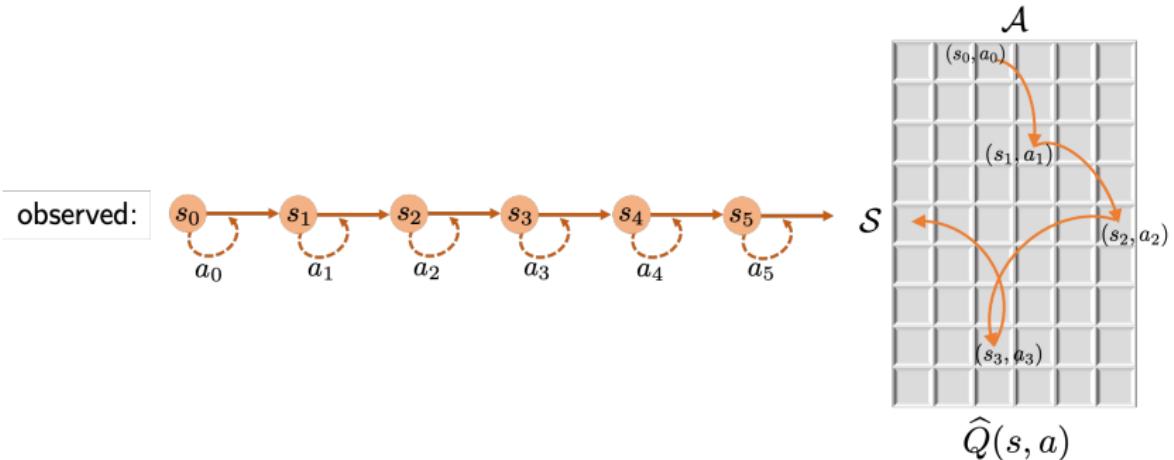
$$\mathcal{T}(Q)(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ \max_{a'} Q(s', a') \right]$$

# Q-learning on Markovian samples



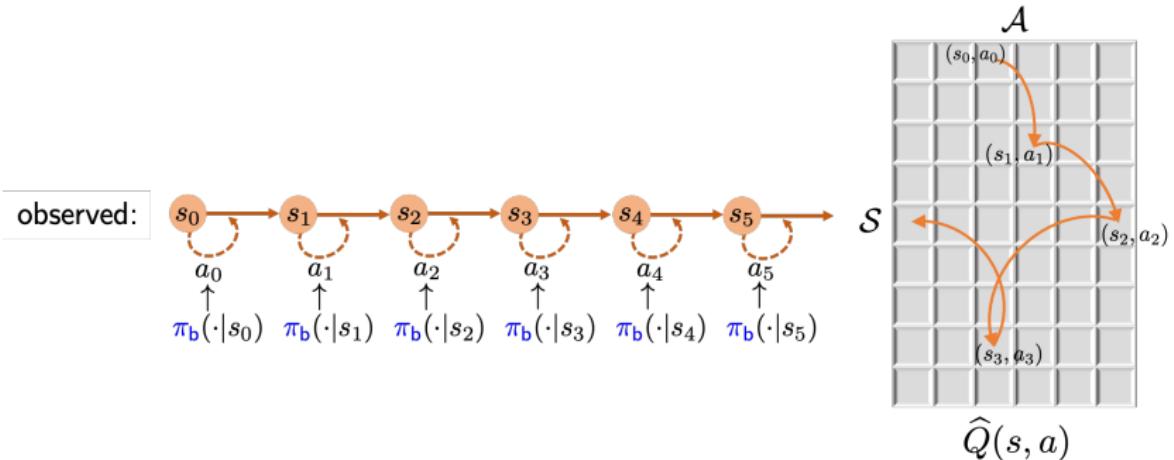
- **asynchronous:** only a single entry is updated each iteration

# Q-learning on Markovian samples



- **asynchronous:** only a single entry is updated each iteration
  - resembles Markov-chain *coordinate descent*

# Q-learning on Markovian samples



- **asynchronous:** only a single entry is updated each iteration
  - resembles Markov-chain *coordinate descent*
- **off-policy:** target policy  $\pi^* \neq$  behavior policy  $\pi_b$

# A highly incomplete list of prior work

---

- Watkins, Dayan '92
- Tsitsiklis '94
- Jaakkola, Jordan, Singh '94
- Szepesvári '98
- Kearns, Singh '99
- Borkar, Meyn '00
- Even-Dar, Mansour '03
- Beck, Srikant '12
- Chi, Zhu, Bubeck, Jordan '18
- Shah, Xie '18
- Lee, He '18
- Wainwright '19
- Chen, Zhang, Doan, Maguluri, Clarke '19
- Yang, Wang '19
- Du, Lee, Mahajan, Wang '20
- Chen, Maguluri, Shakkottai, Shanmugam '20
- Qu, Wierman '20
- Devraj, Meyn '20
- Weng, Gupta, He, Ying, Srikant '20
- ...

*What is sample complexity of (async) Q-learning?*

## Prior art: async Q-learning

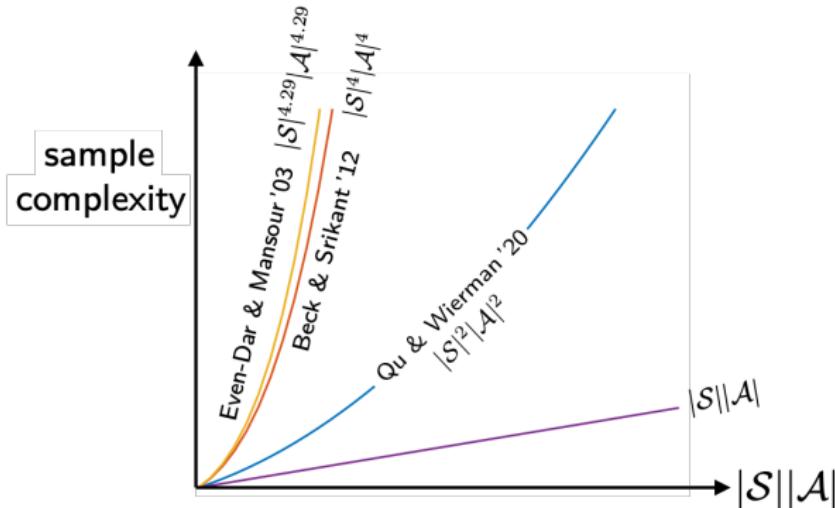
---

**Question:** how many samples are needed to ensure  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$ ?

paper	sample complexity	learning rate
Even-Dar & Mansour '03	$\frac{(t_{\text{cover}})^{\frac{1}{1-\gamma}}}{(1-\gamma)^4 \varepsilon^2}$	linear: $\frac{1}{t}$
Even-Dar & Mansour '03	$\left(\frac{t_{\text{cover}}^{1+3\omega}}{(1-\gamma)^4 \varepsilon^2}\right)^{\frac{1}{\omega}} + \left(\frac{t_{\text{cover}}}{1-\gamma}\right)^{\frac{1}{1-\omega}}$	poly: $\frac{1}{t^\omega}$ , $\omega \in (\frac{1}{2}, 1)$
Beck & Srikant '12	$\frac{t_{\text{cover}}^3  \mathcal{S}   \mathcal{A} }{(1-\gamma)^5 \varepsilon^2}$	constant
Qu & Wierman '20	$\frac{t_{\text{mix}}}{\mu_{\min}^2 (1-\gamma)^5 \varepsilon^2}$	rescaled linear

# Prior art: async Q-learning

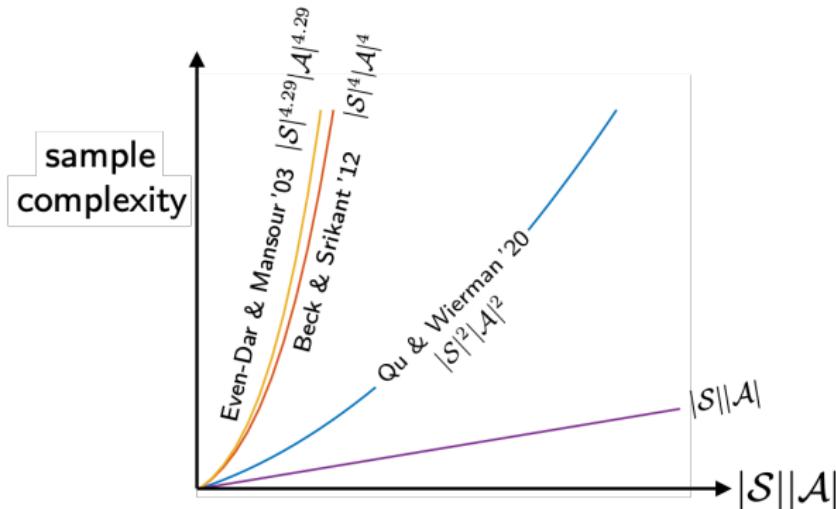
**Question:** how many samples are needed to ensure  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$ ?



if we take  $\mu_{\min} \asymp \frac{1}{|S||\mathcal{A}|}$ ,  $t_{\text{cover}} \asymp \frac{t_{\text{mix}}}{\mu_{\min}}$

# Prior art: async Q-learning

**Question:** how many samples are needed to ensure  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$ ?

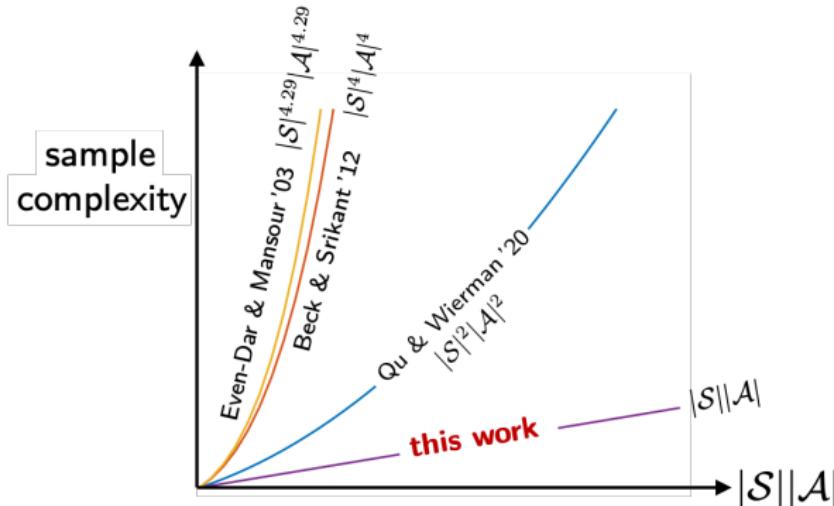


if we take  $\mu_{\min} \asymp \frac{1}{|S||\mathcal{A}|}$ ,  $t_{\text{cover}} \asymp \frac{t_{\text{mix}}}{\mu_{\min}}$

All prior results require sample size of at least  $t_{\text{mix}}|S|^2|\mathcal{A}|^2$ !

# Prior art: async Q-learning

**Question:** how many samples are needed to ensure  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$ ?



if we take  $\mu_{\min} \asymp \frac{1}{|S||\mathcal{A}|}$ ,  $t_{\text{cover}} \asymp \frac{t_{\text{mix}}}{\mu_{\min}}$

All prior results require sample size of at least  $t_{\text{mix}} |S|^2 |\mathcal{A}|^2$ !

## Main result: $\ell_\infty$ -based sample complexity

---

### Theorem 5 (Li, Wei, Chi, Gu, Chen '20)

For any  $0 < \varepsilon \leq \frac{1}{1-\gamma}$ , sample complexity of async Q-learning to yield  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$  is at most (up to some log factor)

$$\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$$

# Main result: $\ell_\infty$ -based sample complexity

## Theorem 5 (Li, Wei, Chi, Gu, Chen '20)

For any  $0 < \varepsilon \leq \frac{1}{1-\gamma}$ , sample complexity of async Q-learning to yield  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$  is at most (up to some log factor)

$$\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$$

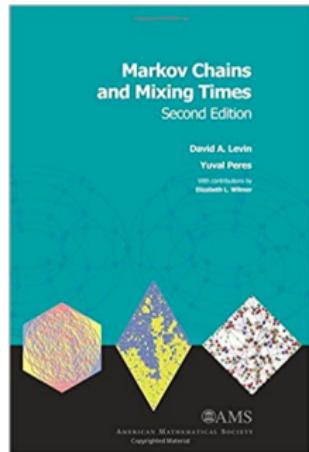
- Improves upon prior art by **at least**  $|\mathcal{S}||\mathcal{A}|$ !

— prior art:  $\frac{t_{\text{mix}}}{\mu_{\min}^2(1-\gamma)^5 \varepsilon^2}$  (Qu & Wierman '20)

# Effect of mixing time on sample complexity

---

$$\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$$

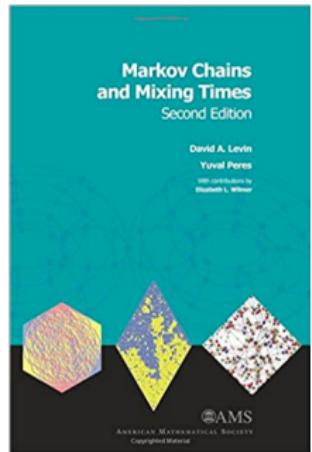


- reflects cost taken to reach steady state
- one-time expense (almost independent of  $\varepsilon$ )
  - it becomes amortized as algorithm runs

# Effect of mixing time on sample complexity

---

$$\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$$



- reflects cost taken to reach steady state
- one-time expense (almost independent of  $\varepsilon$ )
  - it becomes amortized as algorithm runs

— prior art:  $\frac{t_{\text{mix}}}{\mu_{\min}^2(1-\gamma)^5 \varepsilon^2}$  (Qu & Wierman '20)

# Learning rates

---

**Our choice:** constant stepsize  $\eta_t \equiv \min \left\{ \frac{(1-\gamma)^4 \varepsilon^2}{\gamma^2}, \frac{1}{t_{\text{mix}}} \right\}$

- Qu & Wierman '20: rescaled linear  $\eta_t = \frac{\frac{1}{\mu_{\min}(1-\gamma)}}{t + \max\left\{\frac{1}{\mu_{\min}(1-\gamma)}, t_{\text{mix}}\right\}}$
- Beck & Srikant '12: constant  $\eta_t \equiv \underbrace{\frac{(1-\gamma)^4 \varepsilon^2}{|\mathcal{S}||\mathcal{A}|t_{\text{cover}}^2}}_{\text{too conservative}}$
- Even-Dar & Mansour '03: polynomial  $\eta_t = t^{-\omega}$  ( $\omega \in (\frac{1}{2}, 1]$ )

# Minimax lower bound

---

minimax lower bound  
(Azar et al. '13)

$$\frac{1}{\mu_{\min}(1-\gamma)^3 \varepsilon^2}$$

asyn Q-learning  
(ignoring dependency on  $t_{\text{mix}}$ )

$$\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2}$$

# Minimax lower bound

---

minimax lower bound  
(Azar et al. '13)

$$\frac{1}{\mu_{\min}(1-\gamma)^3 \varepsilon^2}$$

asyn Q-learning  
(ignoring dependency on  $t_{\text{mix}}$ )

$$\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2}$$

Can we improve dependency on **discount complexity**  $\frac{1}{1-\gamma}$ ?

# One strategy: variance reduction

---

— inspired by Johnson & Zhang '13, Wainwright '19

## Variance-reduced Q-learning updates

$$Q_t(s_t, a_t) = (1 - \eta)Q_{t-1}(s_t, a_t) + \eta \left( \mathcal{T}_t(Q_{t-1}) \underbrace{- \mathcal{T}_t(\bar{Q}) + \tilde{\mathcal{T}}(\bar{Q})}_{\text{use } \bar{Q} \text{ to help reduce variability}} \right)(s_t, a_t)$$

- $\bar{Q}$ : some reference Q-estimate
- $\tilde{\mathcal{T}}$ : empirical Bellman operator (using a batch of samples)

# Variance-reduced Q-learning

---

— inspired by Johnson & Zhang '13, Sidford et al. '18, Wainwright '19

update variance-reduced

$\bar{Q}$        $Q$ -learning



**for** each epoch

1. update  $\bar{Q}$  and  $\tilde{T}(\bar{Q})$
2. run variance-reduced  $Q$ -learning updates

# Main result: $\ell_\infty$ -based sample complexity

## Theorem 6 (Li, Wei, Chi, Gu, Chen '20)

For any  $0 < \varepsilon \leq 1$ , sample complexity for (async) variance-reduced **Q-learning** to yield  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$  is at most on the order of

$$\frac{1}{\mu_{\min}(1-\gamma)^3 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$$

- more aggressive learning rates:  $\eta_t \equiv \min \left\{ \frac{(1-\gamma)^4(1-\gamma)^2}{\gamma^2}, \frac{1}{t_{\text{mix}}} \right\}$

# Main result: $\ell_\infty$ -based sample complexity

## Theorem 6 (Li, Wei, Chi, Gu, Chen '20)

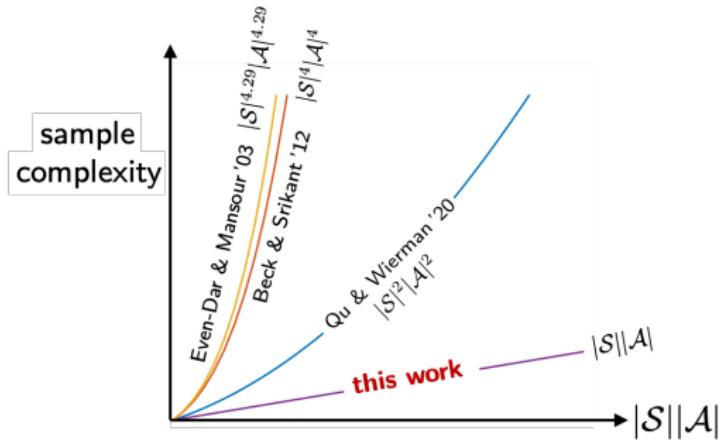
For any  $0 < \varepsilon \leq 1$ , sample complexity for (async) variance-reduced **Q-learning** to yield  $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$  is at most on the order of

$$\frac{1}{\mu_{\min}(1-\gamma)^3 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$$

- more aggressive learning rates:  $\eta_t \equiv \min \left\{ \frac{(1-\gamma)^4(1-\gamma)^2}{\gamma^2}, \frac{1}{t_{\text{mix}}} \right\}$
- minimax-optimal for  $0 < \varepsilon \leq 1$

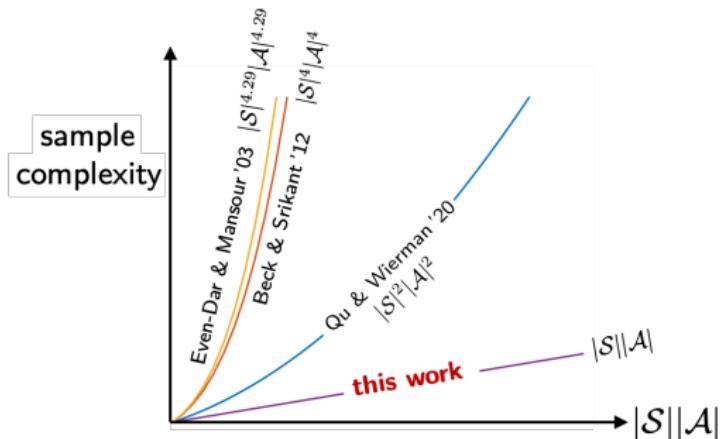
# Summary

Sharpens finite-sample understanding of Q-learning on Markovian data



# Summary

Sharpens finite-sample understanding of Q-learning on Markovian data



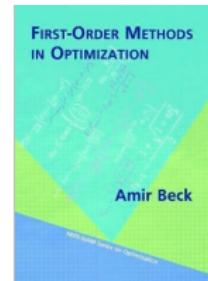
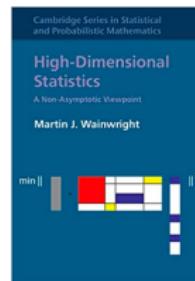
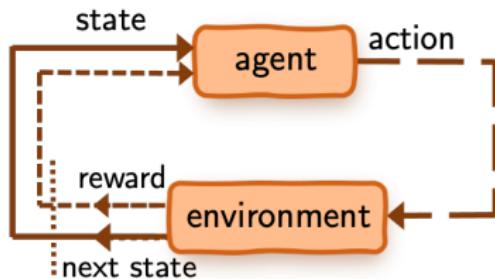
## future directions

- function approximation
- on-policy algorithms like SARSA
- general Markov-chain-based optimization algorithms

# Concluding remarks

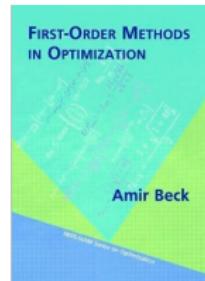
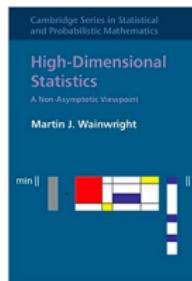
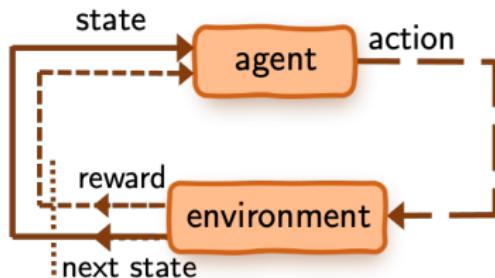
---

Understanding RL requires modern statistics and optimization



# Concluding remarks

Understanding RL requires modern statistics and optimization



## future directions

- beyond tabular settings
- finite-horizon episodic MDPs
- multi-agent RL (e.g. Markov games)
- ...

## Papers:

"Breaking the sample size barrier in model-based reinforcement learning with a generative model," G. Li, Y. Wei, Y. Chi, Y. Gu, Y. Chen, NeurIPS, 2020

"Sample complexity of asynchronous Q-learning: Sharper analysis and variance reduction," G. Li, Y. Wei, Y. Chi, Y. Gu, Y. Chen, NeurIPS 2020

"Fast global convergence of natural policy gradient methods with entropy regularization," S. Cen, C. Cheng, Y. Chen, Y. Wei, Y. Chi, arxiv:2007.06558, 2020