

# Breaking the sample size barrier in reinforcement learning via model-based methods “plug-in”



Yuxin Chen

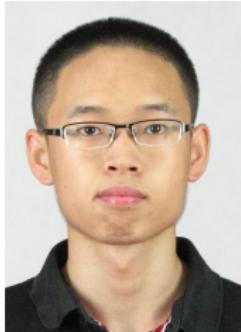
EE, Princeton University

## David Blackwell, 1919–2010: An explorer in mathematics and statistics

Peter J. Bickel<sup>p,1</sup>

Blackwell channel. He also began to work in dynamic programming, which is now called reinforcement learning.<sup>1</sup> In a series of papers, Blackwell gave a rigorous foundation to the theory of dynamic programming, introducing what have become known as Blackwell optimal policies.

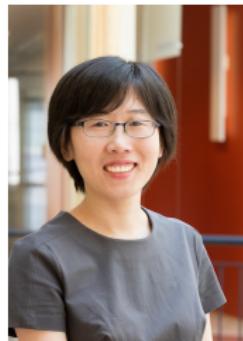




Gen Li  
Tsinghua EE



Yuting Wei  
CMU Statistics

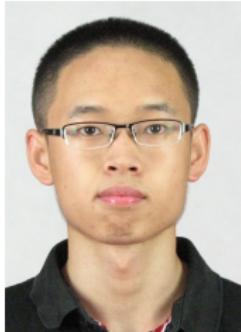


Yuejie Chi  
CMU ECE



Yuantao Gu  
Tsinghua EE

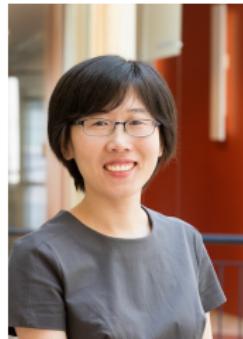
“Breaking the sample size barrier in model-based reinforcement learning with a generative model,” G. Li, Y. Wei, Y. Chi, Y. Gu, Y. Chen, arxiv:2005.12900, 2020



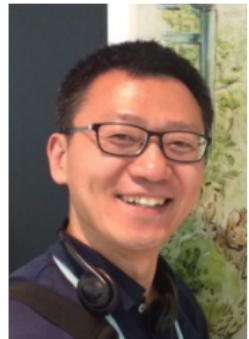
Gen Li  
Tsinghua EE



Yuting Wei  
Berkeley Stat Ph.D.



Yuejie Chi  
CMU ECE

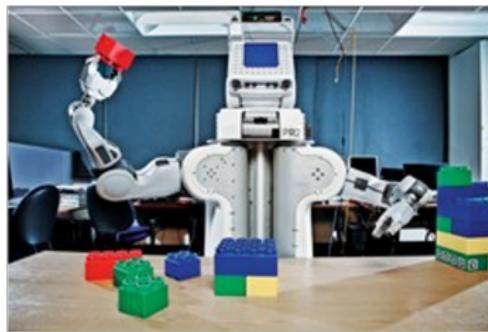


Yuantao Gu  
Tsinghua EE

"Breaking the sample size barrier in model-based reinforcement learning with a generative model," G. Li, Y. Wei, Y. Chi, Y. Gu, Y. Chen, arxiv:2005.12900, 2020

# 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



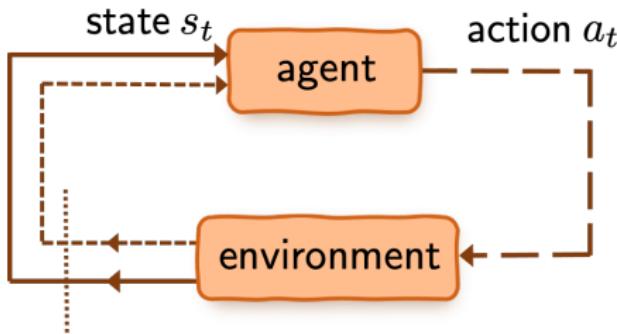
online ads

Calls for design of sample-efficient RL algorithms!

*Background: Markov decision processes*

# Markov decision process (MDP)

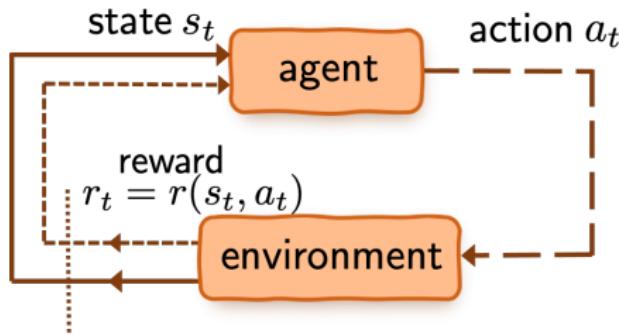
---



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

# Markov decision process (MDP)

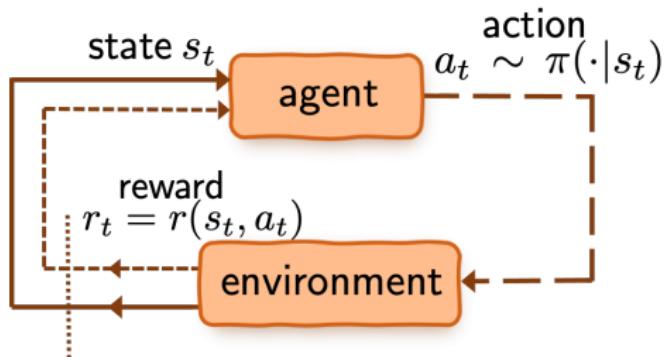
---



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

# Markov decision process (MDP)

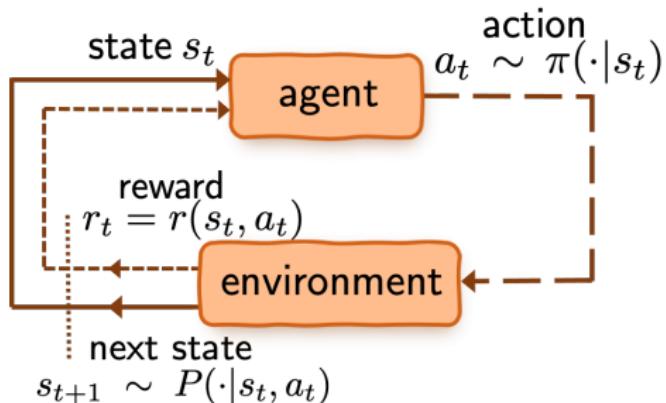
---



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

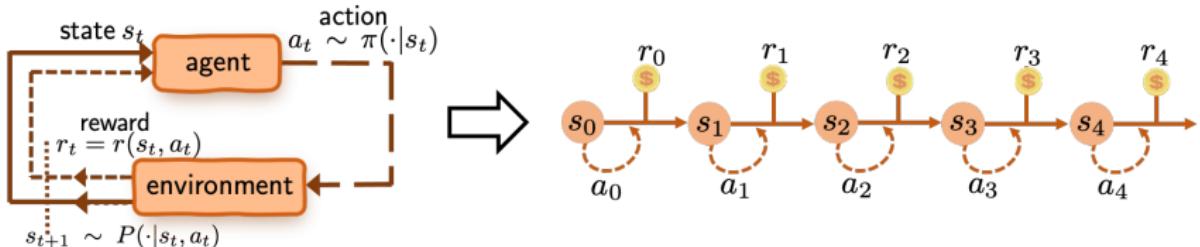
# Markov decision process (MDP)

---



- $\mathcal{S}$ : 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

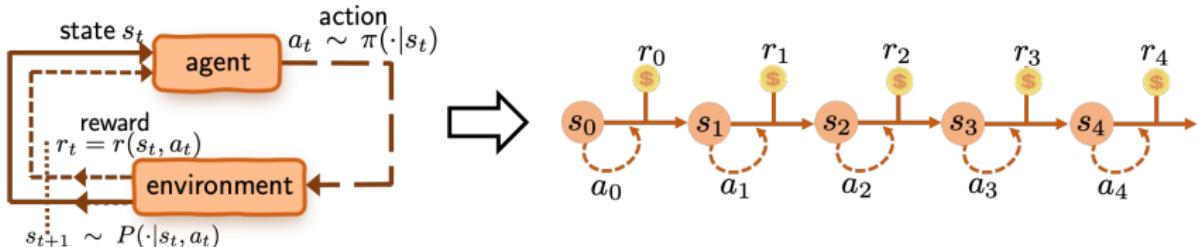


Value of policy  $\pi$ : long-term *discounted* reward

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



# Value function



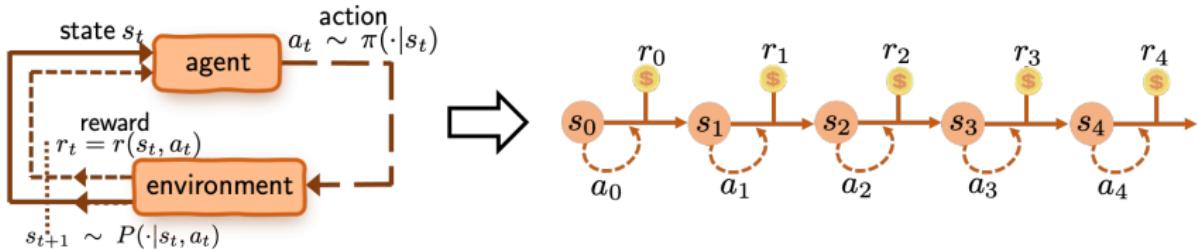
Value of policy  $\pi$ : long-term *discounted* reward

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



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

# Value function



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

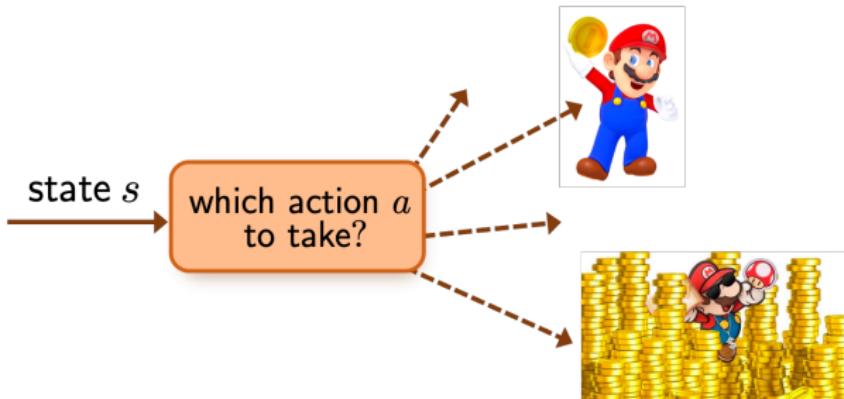
$$\forall s \in \mathcal{S} : V^\pi(s) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right]$$



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

# Optimal policy and optimal values

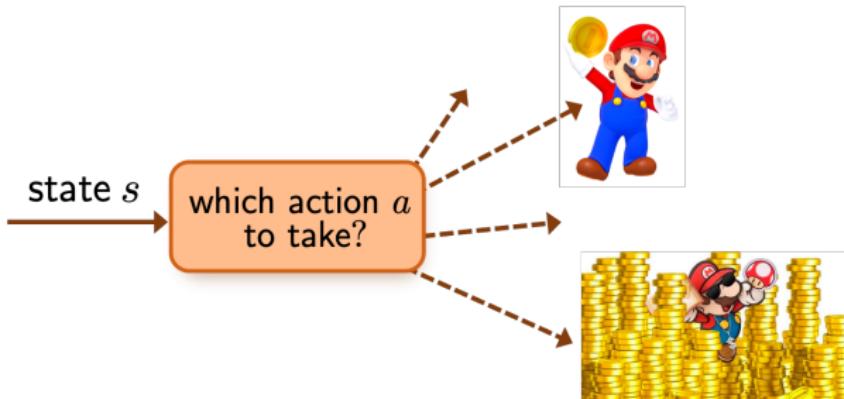
---



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

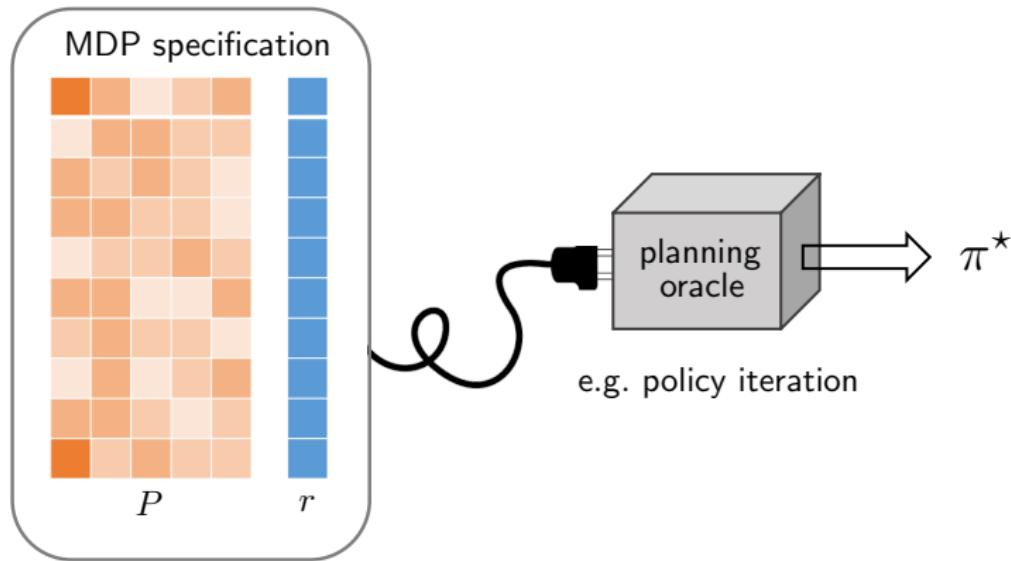
# Optimal policy and optimal values

---



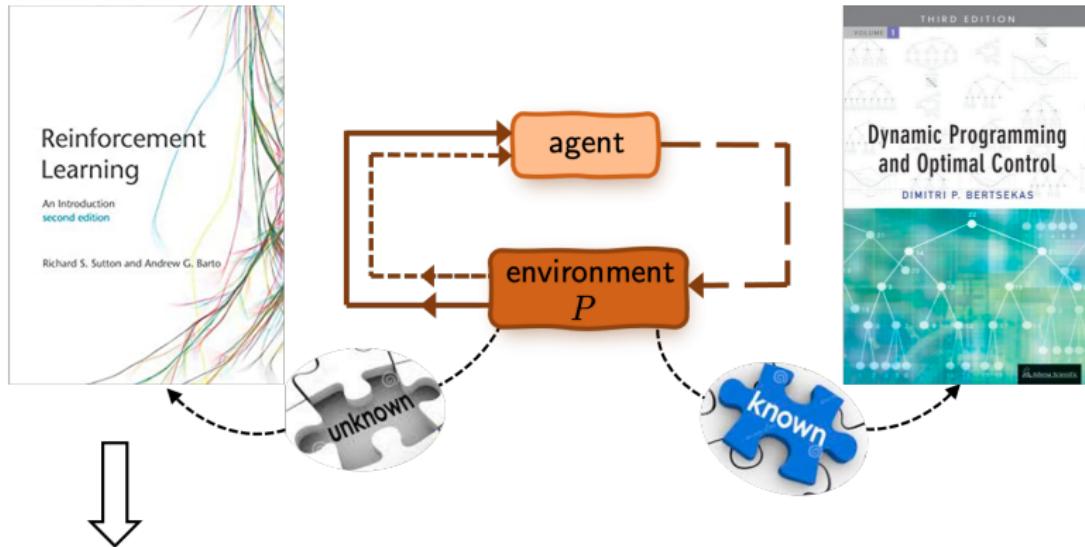
- **Optimal policy**  $\pi^*$ : maximizing the value function
- Optimal values:  $V^* := V^{\pi^*}$

# When the model is known . . .



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

# When the model is unknown ...

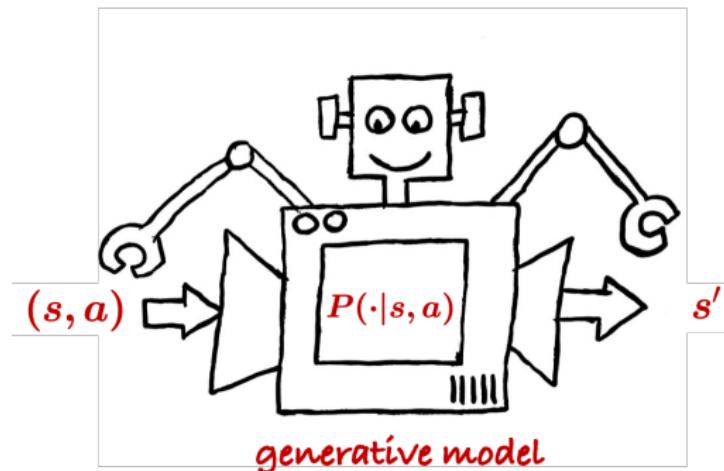


Need to learn optimal policy from samples w/o model specification

# This talk: 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 ?

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

# An incomplete list of prior art

---

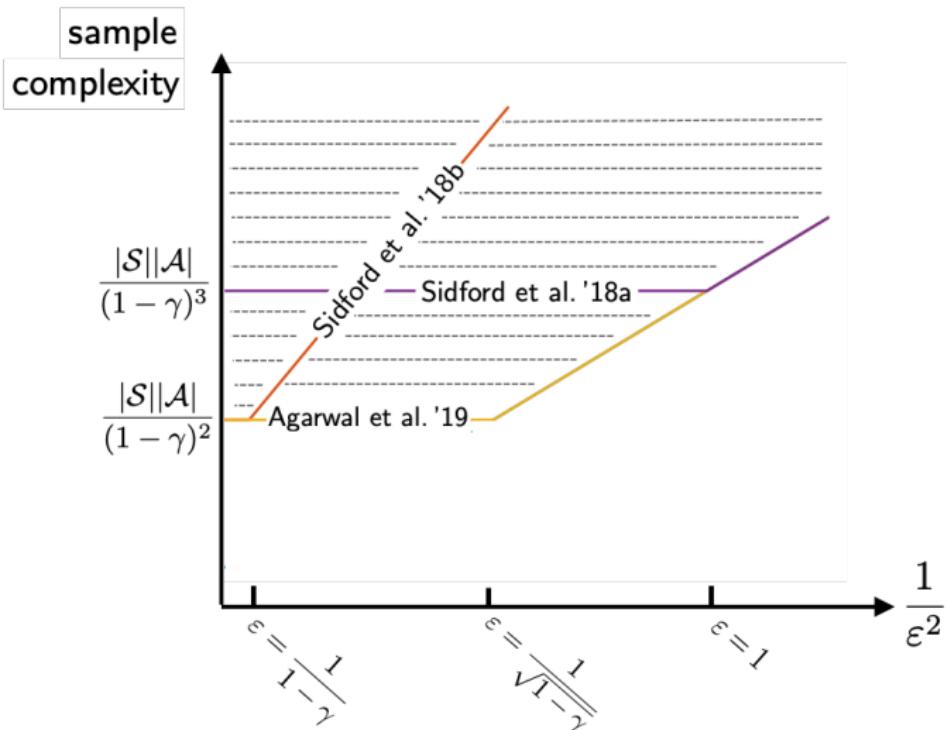
- Kearns & Singh '99
- Kakade '03
- Kearns, Mansour & Ng '02
- Azar, Munos & Kappen '12
- Azar, Munos, Ghavamzadeh & Kappen '13
- Sidford, Wang, Wu, Yang & Ye '18
- Sidford, Wang, Wu & Ye '18
- Wang '17
- Agarwal, Kakade & Yang '19
- Wainwright '19a
- Wainwright '19b
- Pananjady & Wainwright '20
- Yang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

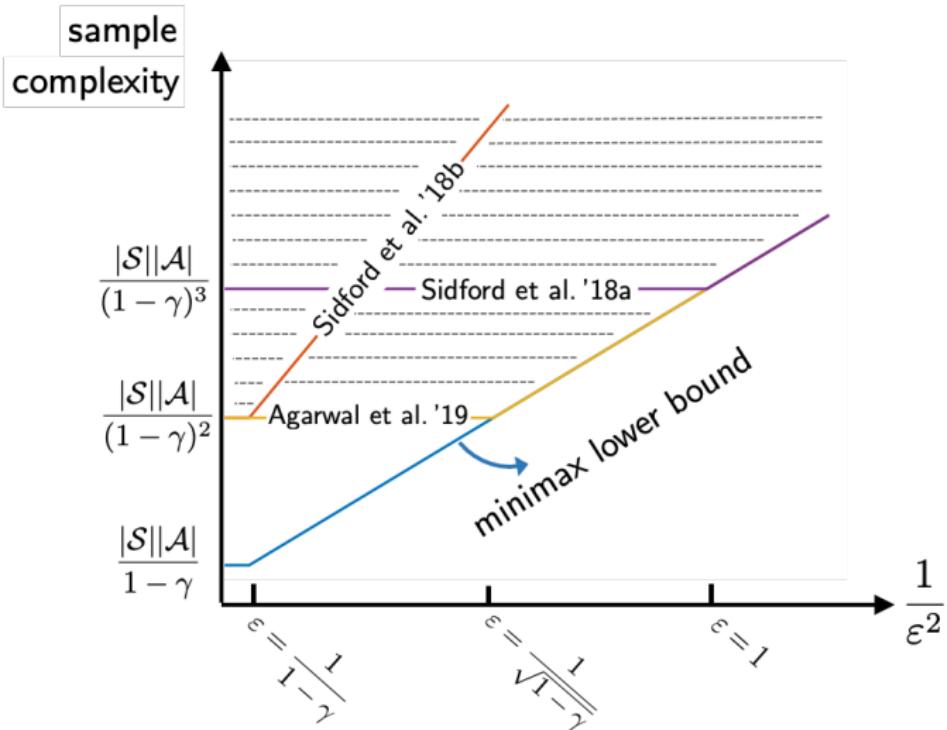
## An even shorter list of prior art

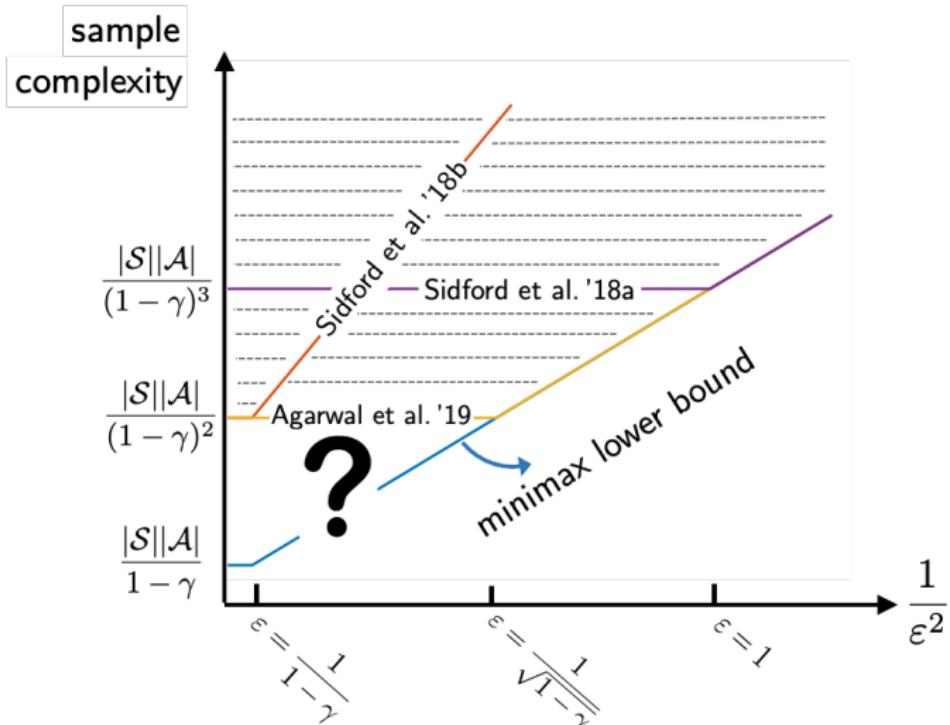
---

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)





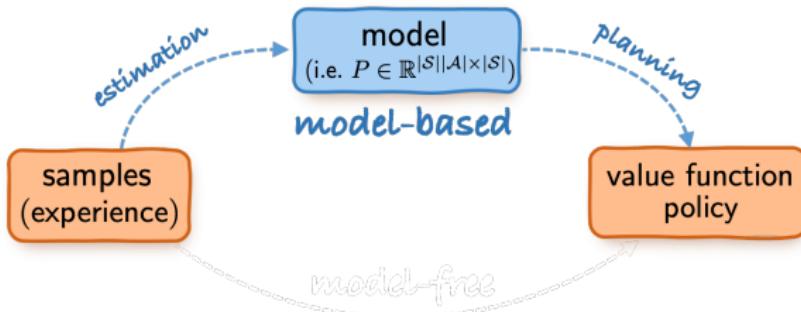


All prior theory requires sample size  $> \underbrace{\frac{|S||\mathcal{A}|}{(1 - \gamma)^2}}_{\text{sample size barrier}}$

*Is it possible to close the gap?*

# Two approaches

---

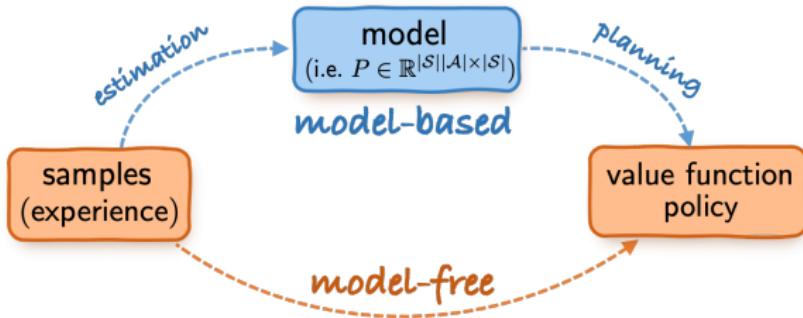


## Model-based approach (“plug-in”)

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

# Two approaches

---



## Model-based approach (“plug-in”)

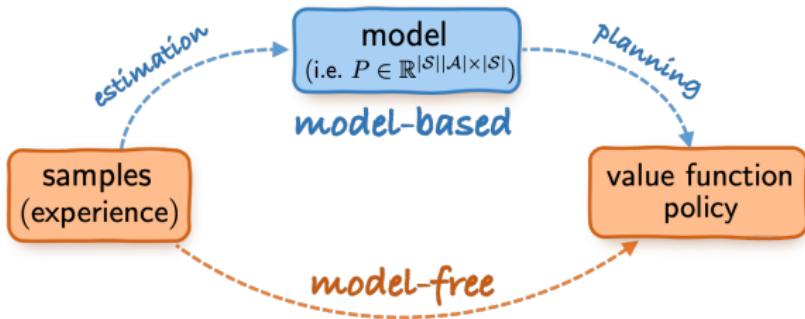
1. build an empirical estimate  $\hat{P}$  for  $P$
2. planning based on the empirical  $\hat{P}$

## Model-free approach (e.g. Q-learning, SARSA)

— learning w/o estimating the model explicitly

# Two approaches

---



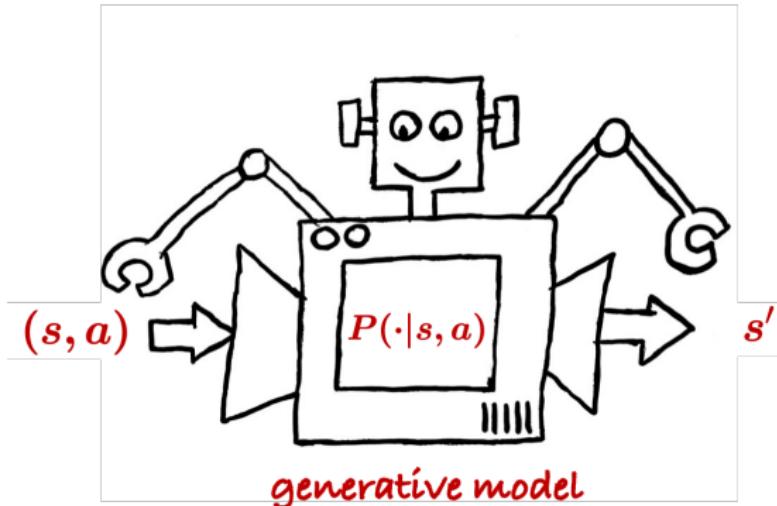
## Model-based approach (“plug-in”)

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

## Model-free approach (e.g. Q-learning, SARSA)

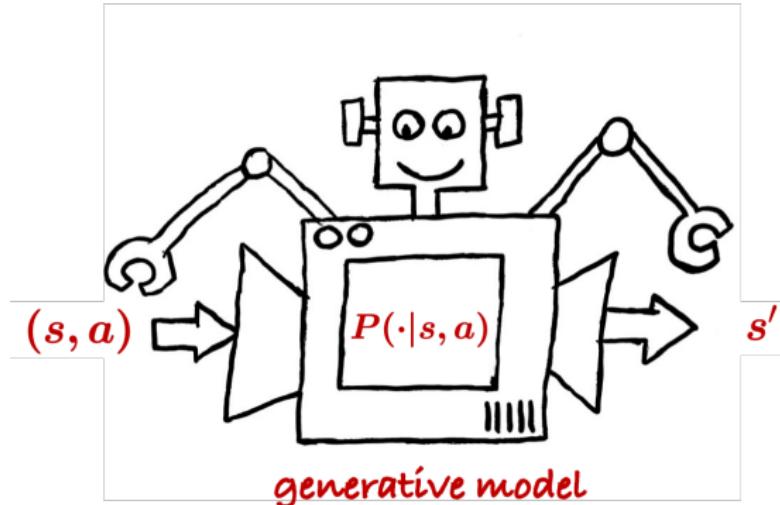
— learning w/o estimating the model explicitly

# Model estimation



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

# Model estimation

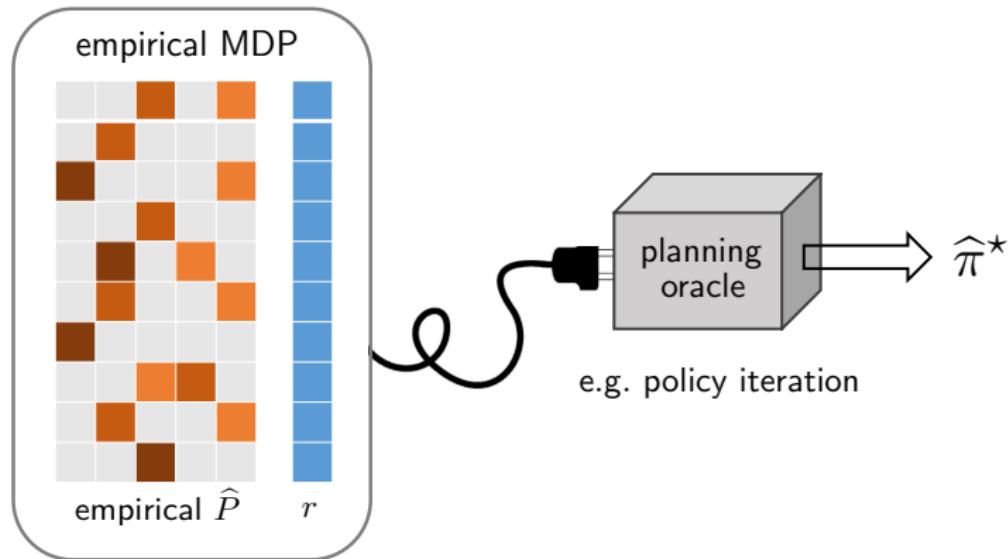


**Sampling:** for each  $(s, a)$ , collect  $N$  ind. samples  $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$

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

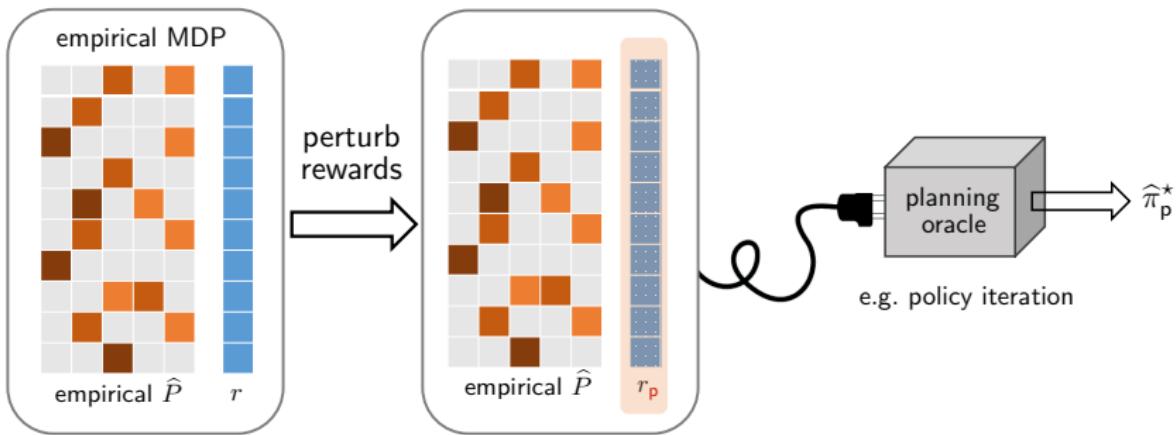
# Model-based (plug-in) estimator

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



Run planning algorithms based on the *empirical* MDP

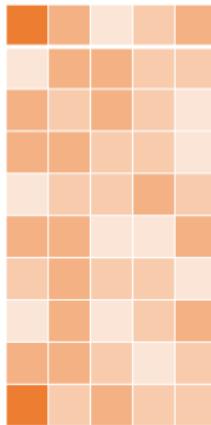
# Our method: plug-in estimator + perturbation



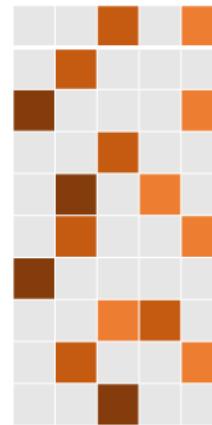
Planning based on the *empirical* MDP with *slightly perturbed rewards*

# Challenges in the sample-starved regime

---



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

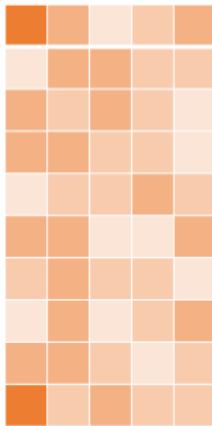


empirical estimate:  
 $\hat{P}$

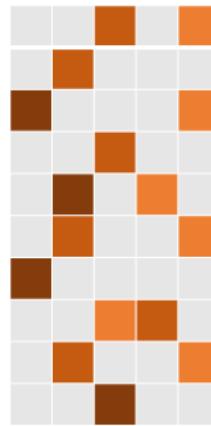
- Can't recover  $P$  faithfully if sample size  $\ll |\mathcal{S}|^2|\mathcal{A}|$ !

# 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)

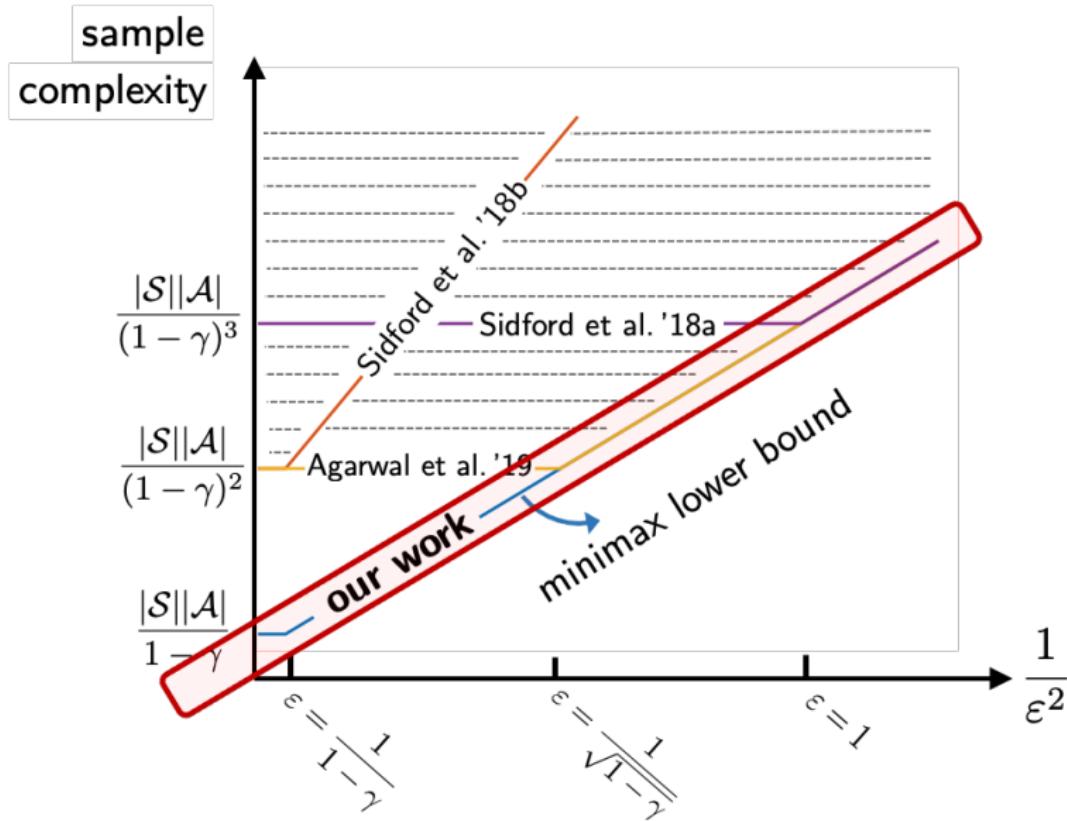
For any  $0 < \varepsilon \leq \frac{1}{1-\gamma}$ , the optimal policy  $\widehat{\pi}_p^*$  of the perturbed empirical MDP achieves

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

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:

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

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

# Step 1: improved theory for policy evaluation

---

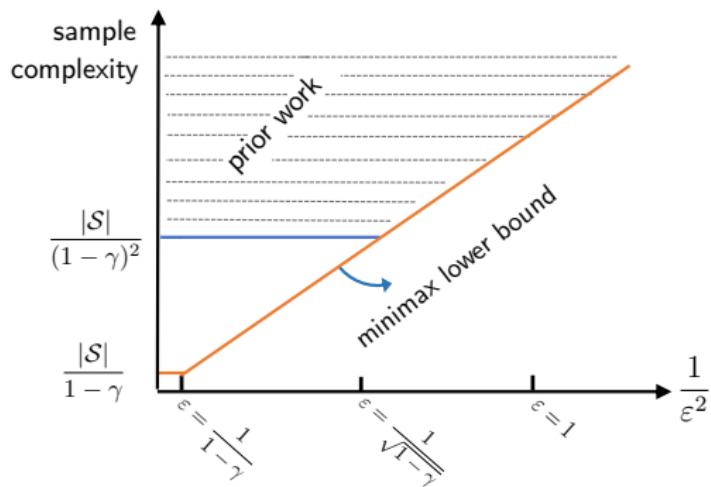
## Model-based policy evaluation:

- 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 $\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}$

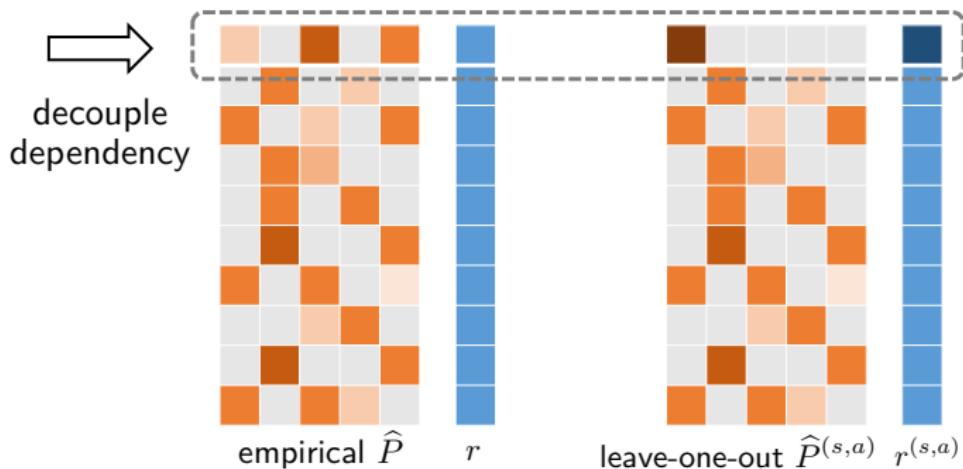
---

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 ...

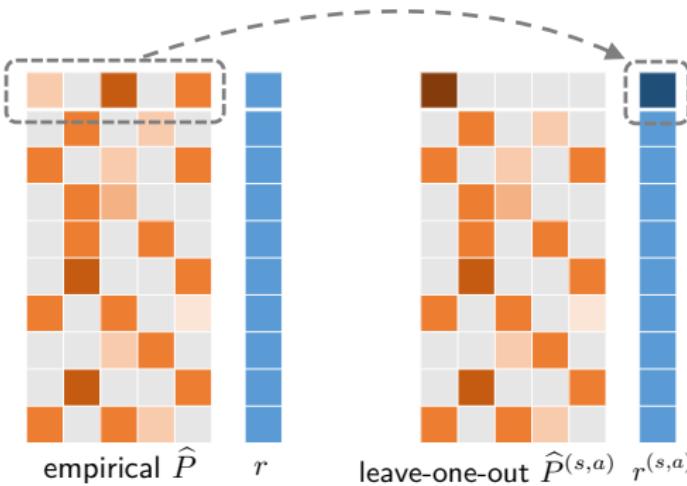
## Key idea 2: leave-one-out analysis

---

- 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

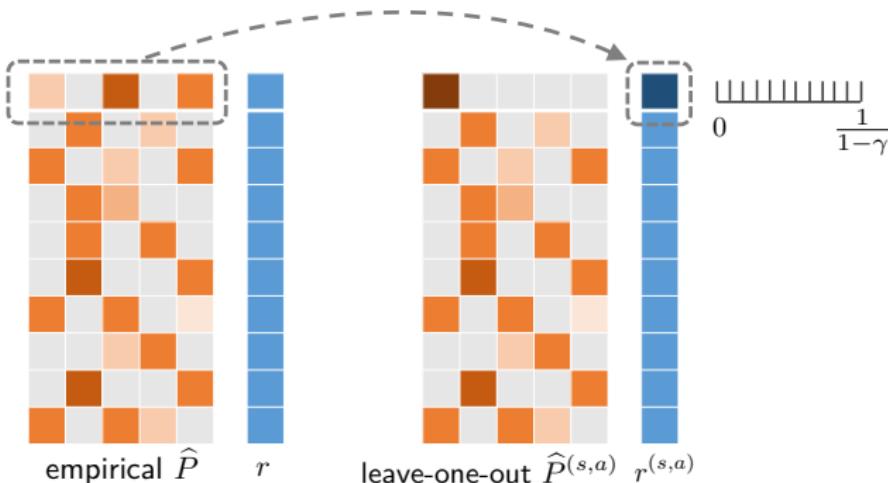
## Key idea 2: leave-one-out analysis

---



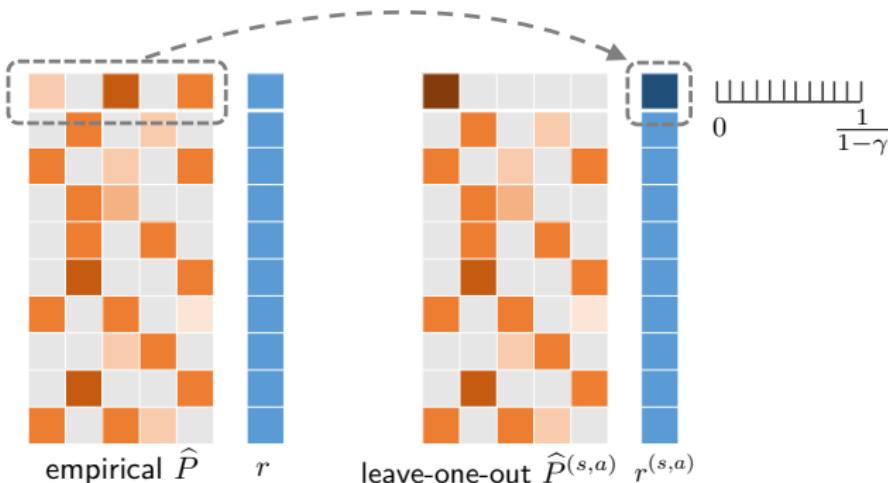
1. embed all randomness from  $\hat{P}_{s,a}$  into a single scalar (i.e.  $r_{s,a}^{(s,a)}$ )

## 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

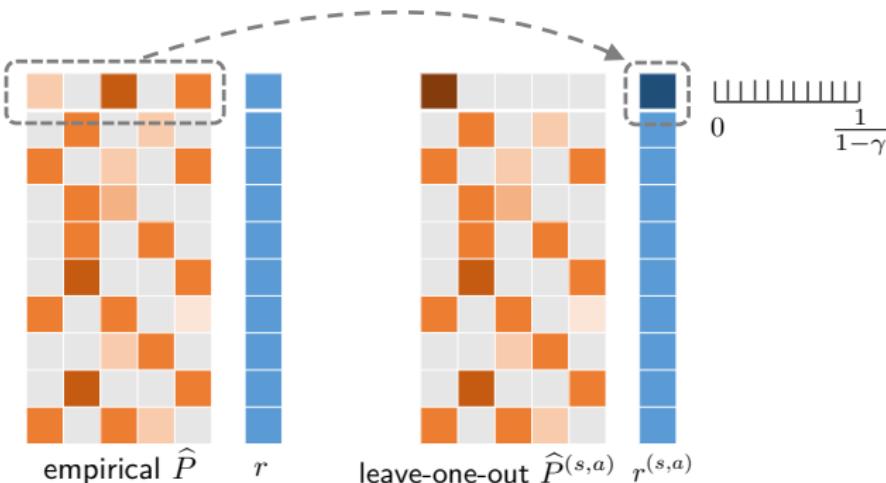
## Key idea 2: leave-one-out analysis



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?

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

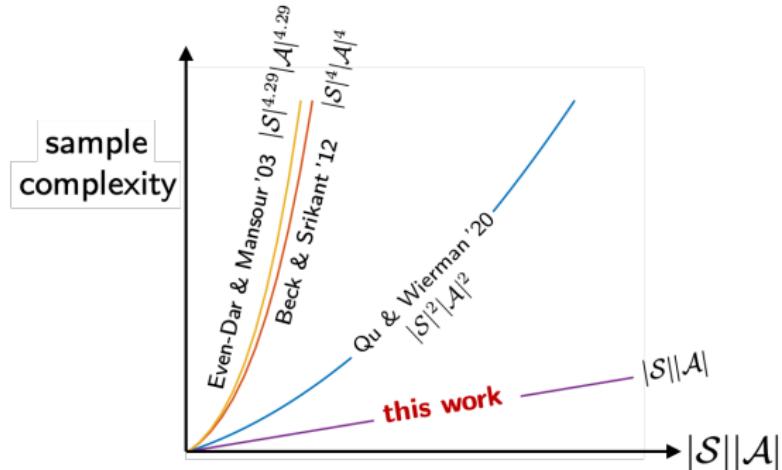
- **Solution:** slightly perturb rewards  $r$   $\implies \hat{\pi}_p^*$

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



## Other stories: sharpened analysis of Q-learning

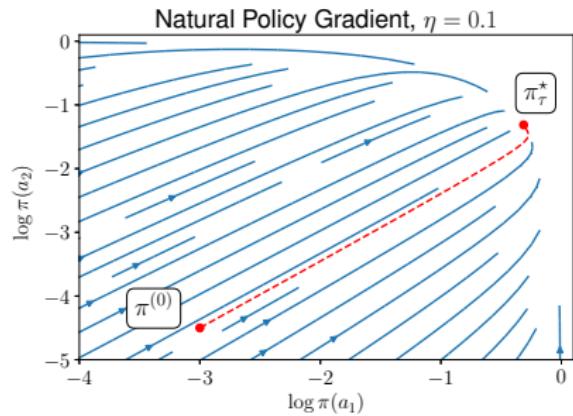
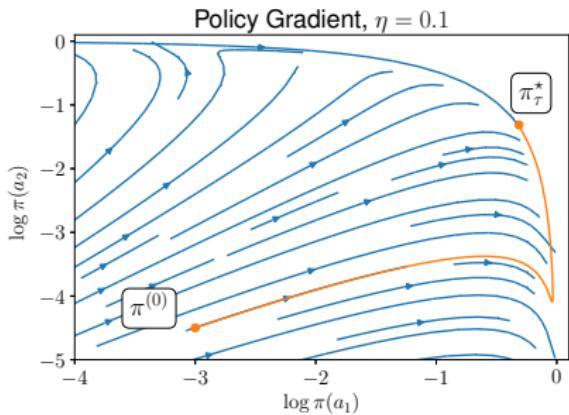
Improves existing sample complexity guarantees for asynchronous Q-learning by at least a factor of  $|\mathcal{S}||\mathcal{A}|!$



"Sample Complexity of Asynchronous Q-Learning: Sharper Analysis and Variance Reduction," G. Li, Y. Wei, Y. Chi, Y. Gu, Y. Chen, NeurIPS 2020

# Other stories: efficiency of natural policy gradient

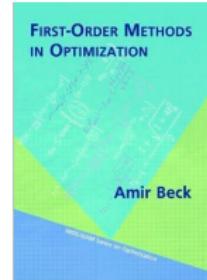
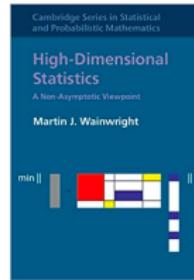
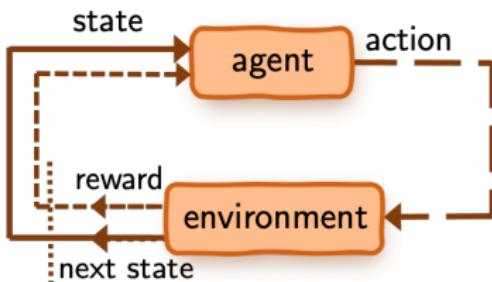
NPG method with entropy regularization converges linearly!



“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

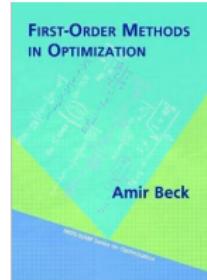
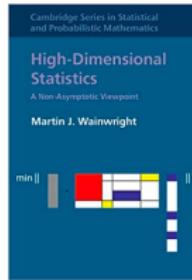
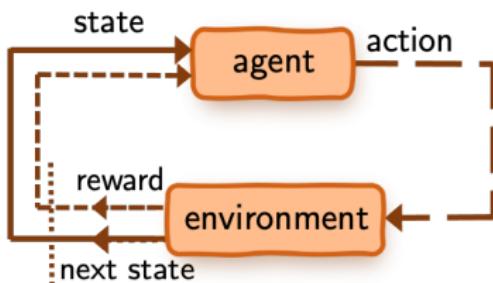
# Concluding remarks

Understanding RL requires modern statistics and optimization



# Concluding remarks

Understanding RL requires modern statistics and optimization



## future directions

- beyond the tabular settings
- finite-horizon episodic MDPs
- Markov games

"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