

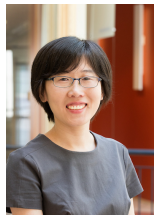
Non-Asymptotic Analysis for Reinforcement Learning



Yuting Wei
UPenn



Yuxin Chen
UPenn



Yuejie Chi
CMU

SIGMETRICS Tutorial, June 2023

Non-asymptotic Analysis for Reinforcement Learning (Part 1)

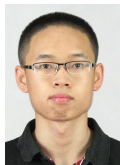


Yuting Wei

Statistics & Data Science, Wharton
University of Pennsylvania

SIGMETRICS, June 2023

Our wonderful collaborators



Gen Li

UPenn → CUHK



Shicong Cen

CMU



Chen Cheng

Stanford



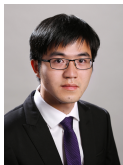
Laixi Shi

CMU → Caltech



Yuling Yan

Princeton → MIT



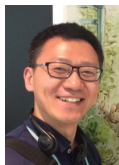
Changxiao Cai

UPenn → UMich



Wenhao Zhan

Princeton



Yuantao Gu

Tsinghua



Jason Lee

Princeton



Jianqing Fan

Princeton

Recent successes in reinforcement learning (RL)



RL holds great promise in the next era of artificial intelligence.

Recap: Supervised learning

Given i.i.d training data, the goal is to make prediction on unseen data:

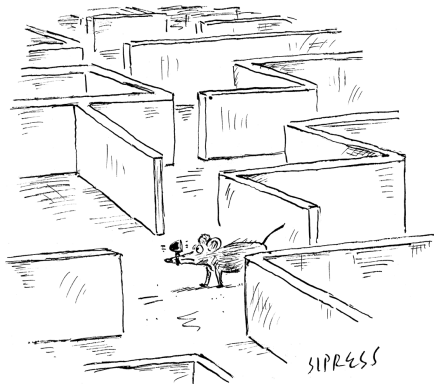


— pic from internet

Reinforcement learning (RL)

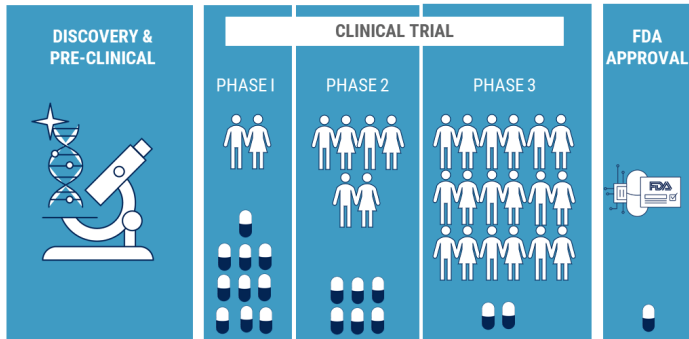
In RL, an agent learns by interacting with an environment.

- no training data
- trial-and-error
- maximize total rewards
- delayed reward



“Recalculating ... recalculating ...”

Sample efficiency

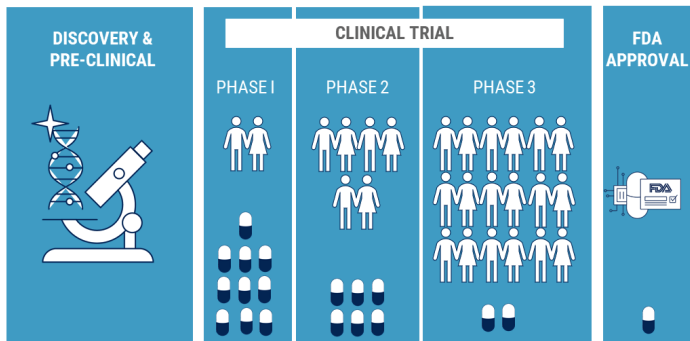


Source: cbinsights.com

CBINSIGHTS

- prohibitively large state & action space
- collecting data samples can be expensive or time-consuming

Sample efficiency



Source: cbinsights.com

CBINSIGHTS

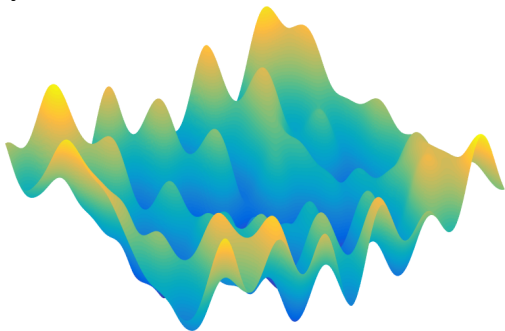
- prohibitively large state & action space
- collecting data samples can be expensive or time-consuming

Challenge: design sample-efficient RL algorithms

Computational efficiency

Running RL algorithms might take a long time ...

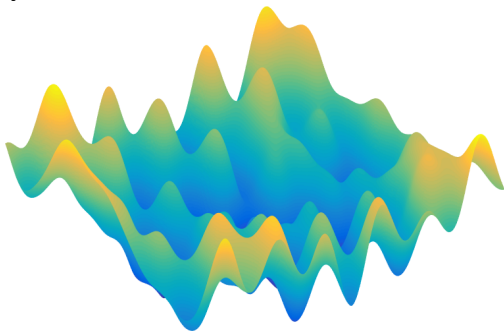
- enormous state-action space
- nonconvexity



Computational efficiency

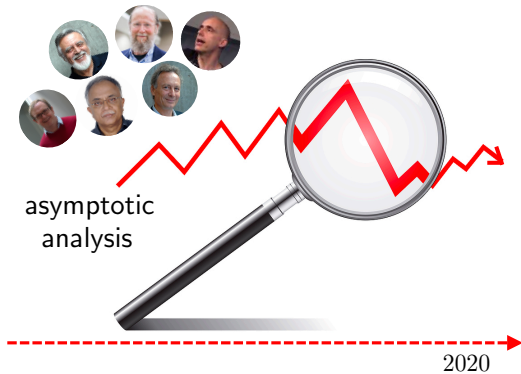
Running RL algorithms might take a long time ...

- enormous state-action space
- nonconvexity

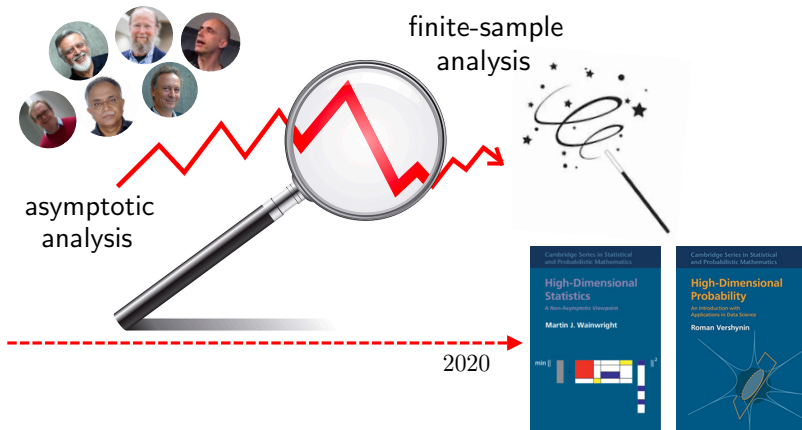


Challenge: design computationally efficient RL algorithms

Theoretical foundation of RL

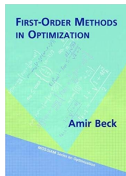
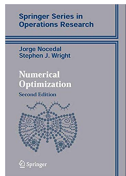


Theoretical foundation of RL

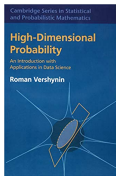
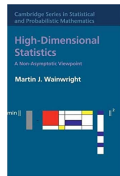


Understanding sample efficiency of RL requires a modern suite of non-asymptotic analysis tools

This tutorial



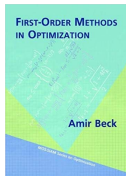
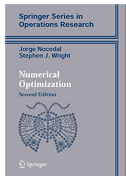
(large-scale) optimization



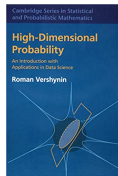
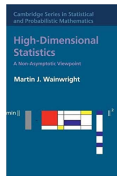
(high-dimensional) statistics

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

This tutorial



(large-scale) optimization



(high-dimensional) statistics

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

Part 1. **basics, and model-based RL**

Part 2. **value-based RL**

Part 3. **policy optimization**

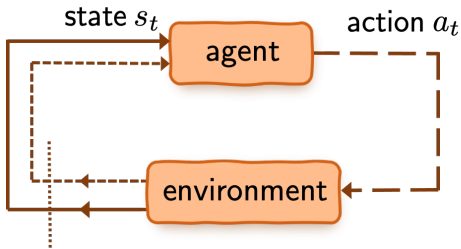
We will illustrate these approaches for learning standard, robust, and multi-agent RL with simulator/online/offline data.

Outline (Part 1)

- Basics: Markov decision processes
- Basic dynamic programming algorithms
- Model-based RL (“plug-in” approach)

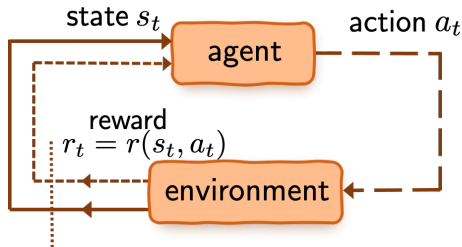
Basics: Markov decision processes

Markov decision process (MDP)



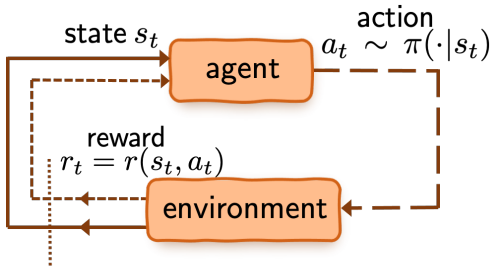
- \mathcal{S} : state space
- \mathcal{A} : action space

Markov decision process (MDP)



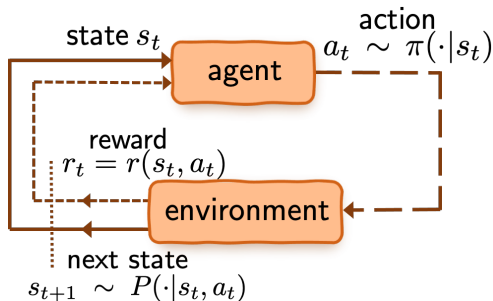
- \mathcal{S} : state space
- \mathcal{A} : action space
- $r(s, a) \in [0, 1]$: immediate reward

Infinite-horizon Markov decision process



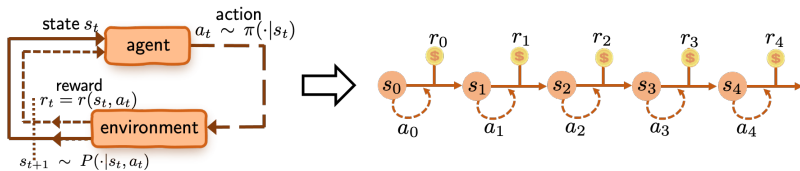
- \mathcal{S} : state space
- \mathcal{A} : action space
- $r(s, a) \in [0, 1]$: immediate reward
- $\pi(\cdot | s)$: policy (or action selection rule)

Infinite-horizon Markov decision process



- \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

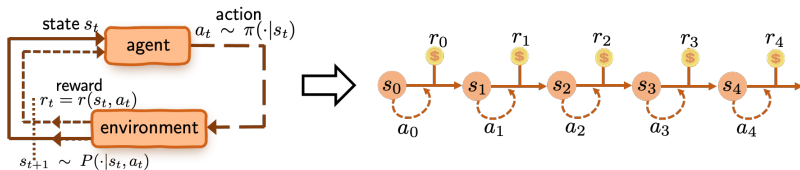
Value function



Value of policy π : cumulative **discounted** reward

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

Value function

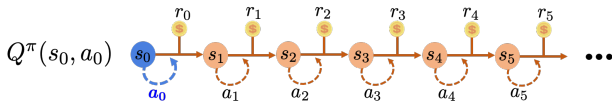


Value of policy π : cumulative **discounted** reward

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

- $\gamma \in [0, 1)$: discount factor
 - ▶ take $\gamma \rightarrow 1$ to approximate **long-horizon** MDPs
 - ▶ **effective horizon**: $\frac{1}{1-\gamma}$

Q-function (action-value function)

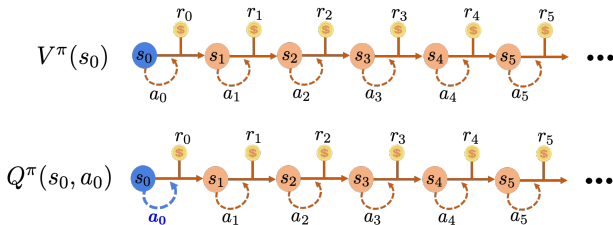


Q-function of policy π :

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

- (~~a_0~~ , $s_1, a_1, s_2, a_2, \dots$): induced by policy π

Q-function (action-value function)

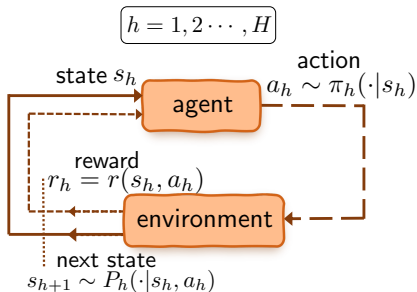


Q-function of policy π :

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

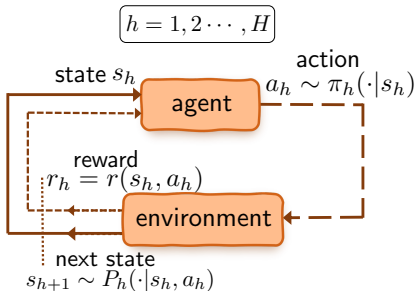
- (~~a₀~~, s₁, a₁, s₂, a₂, ...): induced by policy π

Finite-horizon MDPs



- H : horizon length
- \mathcal{S} : state space with size S
- \mathcal{A} : action space with size A
- $r_h(s_h, a_h) \in [0, 1]$: immediate reward in step h
- $\pi = \{\pi_h\}_{h=1}^H$: policy (or action selection rule)
- $P_h(\cdot | s, a)$: transition probabilities in step h

Finite-horizon MDPs

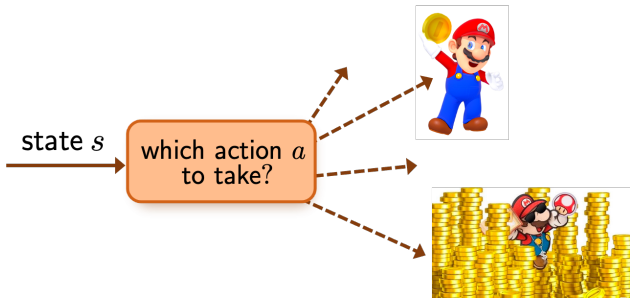


value function: $V_h^\pi(s) := \mathbb{E} \left[\sum_{t=h}^H r_t(s_t, a_t) \mid s_h = s \right]$

Q-function: $Q_h^\pi(s, a) := \mathbb{E} \left[\sum_{t=h}^H r_t(s_t, a_t) \mid s_h = s, a_h = a \right]$



Optimal policy and optimal value



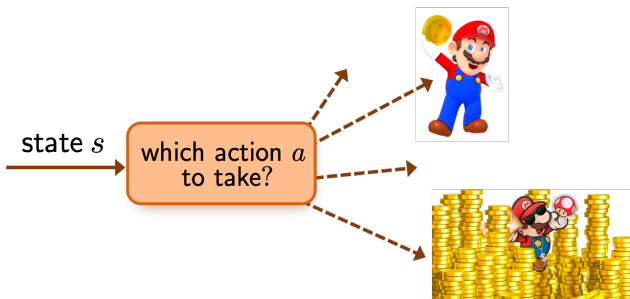
optimal policy π^* : maximizing value function $\max_{\pi} V^{\pi}$

Proposition (Puterman'94)

For infinite horizon discounted MDP, there always exists a deterministic policy π^ , such that*

$$V^{\pi^*}(s) \geq V^{\pi}(s), \quad \forall s, \text{ and } \pi.$$

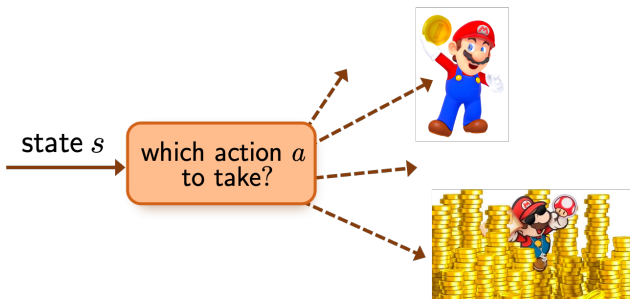
Optimal policy and optimal value



optimal policy π^* : maximizing value function $\max_{\pi} V^{\pi}$

- optimal value / Q function: $V^* := V^{\pi^*}$, $Q^* := Q^{\pi^*}$

Optimal policy and optimal value



optimal policy π^* : maximizing value function $\max_{\pi} V^{\pi}$

- optimal value / Q function: $V^* := V^{\pi^*}$, $Q^* := Q^{\pi^*}$
- How to find this π^* ?

**Basic dynamic programming algorithms
when MDP specification is **known****

Policy evaluation: Given MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, r, P, \gamma)$ and policy $\pi : \mathcal{S} \mapsto \mathcal{A}$, how good is π ? (i.e., how to compute $V^\pi(s)$, $\forall s$?)

Policy evaluation: Given MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, r, P, \gamma)$ and policy $\pi : \mathcal{S} \mapsto \mathcal{A}$, how good is π ? (i.e., how to compute $V^\pi(s)$, $\forall s$?)

Possible scheme:

- execute policy evaluation for each π
- find the optimal one

Policy evaluation: Bellman's consistency equation

- V^π / Q^π : value / action-value function under policy π

Policy evaluation: Bellman's consistency equation

- V^π / Q^π : value / action-value function under policy π

Bellman's consistency equation

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} [Q^\pi(s, a)]$$
$$Q^\pi(s, a) = \underbrace{r(s, a)}_{\text{immediate reward}} + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\underbrace{V^\pi(s')}_{\text{next state's value}} \right]$$



Richard Bellman

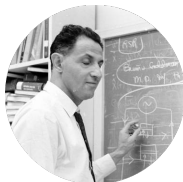
Policy evaluation: Bellman's consistency equation

- V^π / Q^π : value / action-value function under policy π

Bellman's consistency equation

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} [Q^\pi(s, a)]$$
$$Q^\pi(s, a) = \underbrace{r(s, a)}_{\text{immediate reward}} + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\underbrace{V^\pi(s')}_{\text{next state's value}} \right]$$

- one-step look-ahead



Richard Bellman

Policy evaluation: Bellman's consistency equation

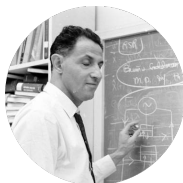
- V^π / Q^π : value / action-value function under policy π

Bellman's consistency equation

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} [Q^\pi(s, a)]$$
$$Q^\pi(s, a) = \underbrace{r(s, a)}_{\text{immediate reward}} + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\underbrace{V^\pi(s')}_{\text{next state's value}} \right]$$

- one-step look-ahead
- let P^π be the state-action transition matrix induced by π :

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



Richard Bellman

Optimal policy π^* : Bellman's 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

Optimal policy π^* : Bellman's 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

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

γ -contraction of Bellman operator:

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



Richard Bellman

Two dynamic programming algorithms

Value iteration (VI)

For $t = 0, 1, \dots$,

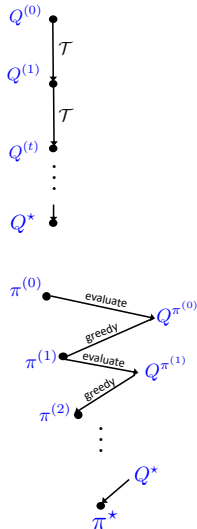
$$Q^{(t+1)} = \mathcal{T}(Q^{(t)})$$

Policy iteration (PI)

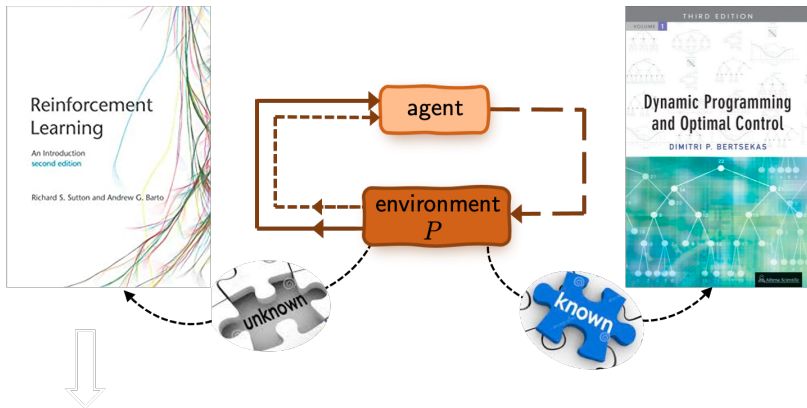
For $t = 0, 1, \dots$,

policy evaluation: $Q^{(t)} = Q^{\pi^{(t)}}$

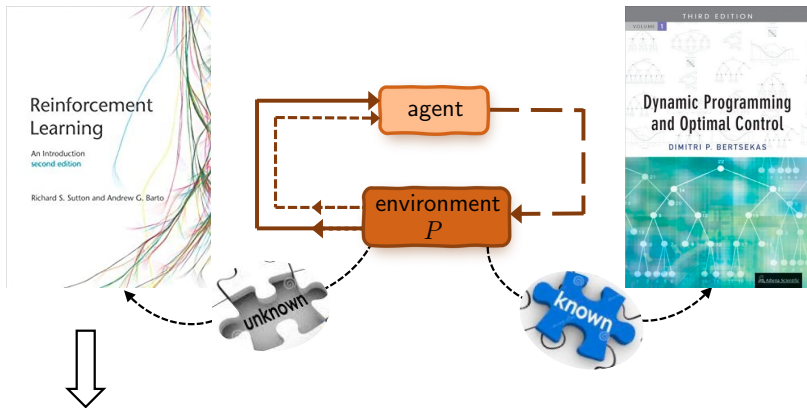
policy improvement: $\pi^{(t+1)}(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q^{(t)}(s, a)$



When the model is unknown ...

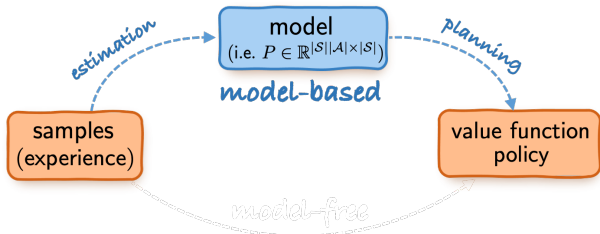


When the model is unknown ...



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

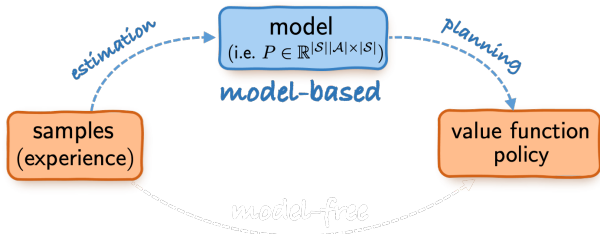
Three approaches



Model-based approach (“plug-in”)

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

Three approaches



Model-based approach (“plug-in”)

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

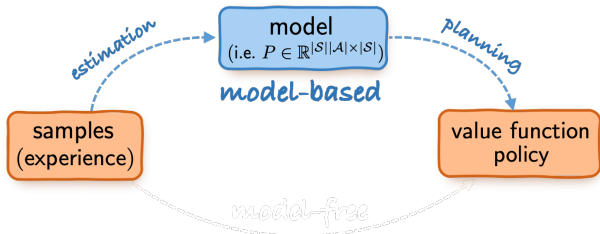
Tutorial Part 2: Value-based approach

— learning w/o estimating the model explicitly

Tutorial Part 3: Policy-based approach

— optimization in the space of policies

Three approaches



Model-based approach (“plug-in”)

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

Tutorial Part 2: Value-based approach

— learning w/o estimating the model explicitly

Tutorial Part 3: Policy-based approach

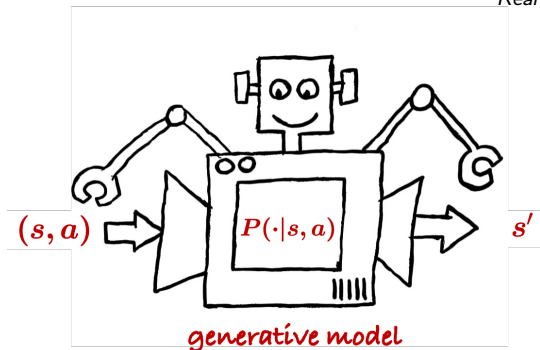
— optimization in the space of policies

Model-based RL (a “plug-in” approach)

1. Sampling from a generative model (simulator)
2. Offline RL / batch RL
3. Robust RL

A generative model / simulator

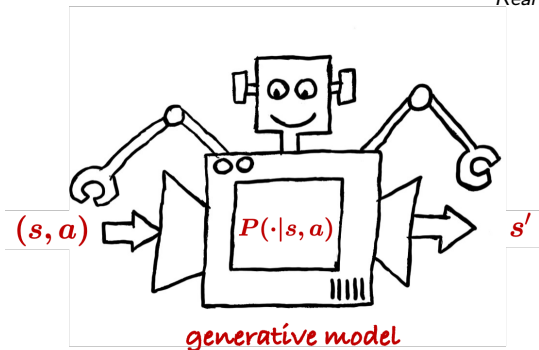
— Kearns and Singh, 1999



- **sampling:** for each (s, a) , collect N samples $\{(s, a, s'_i)\}_{1 \leq i \leq N}$

A generative model / simulator

— Kearns and Singh, 1999



- **sampling:** for each (s, a) , collect N samples $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$
- construct $\hat{\pi}$ based on samples (in total $|\mathcal{S}||\mathcal{A}| \times N$)

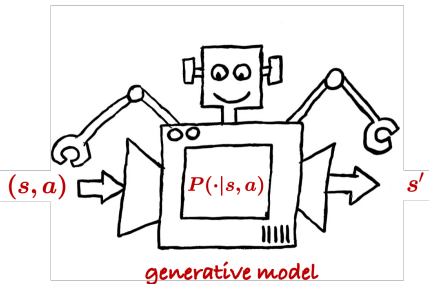
l_∞ -**sample complexity**: how many samples are required to learn an ε -optimal policy?

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

An incomplete list of works

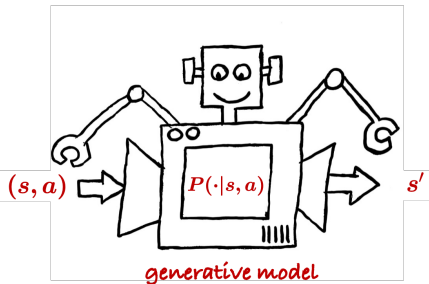
- Kearns and Singh, 1999
- Kakade, 2003
- Kearns 3t al., 2002
- Azar et al., 2012
- **Azar et al., 2013**
- Sidford et al, 2018a, 2018b
- Wang, 2019
- **Agarwal et al, 2019**
- Wainwright, 2019a, 2019b
- Pananjady and Wainwright, 2019
- Yang and Wang, 2019
- Khamaru, 2020
- Mou et al., 2020
- **Li et al., 2020**
- Cui and Yang, 2021
- ...

Model estimation



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

Model estimation



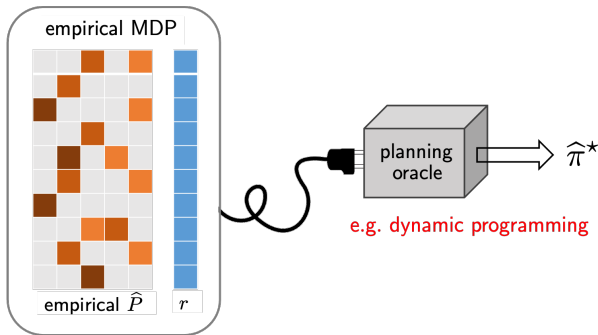
Sampling: for each (s, a) , collect N ind. samples $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$

Empirical estimates:

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

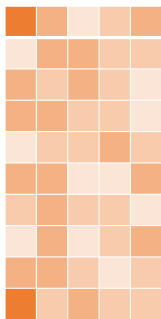
Empirical MDP + planning

— Azar et al., 2013, Agarwal et al., 2019

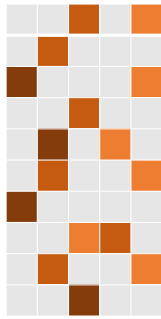


Find policy based on the empirical MDP (*empirical maximizer*)
using, e.g., policy iteration (\hat{P}, r)

Challenges in the sample-starved regime



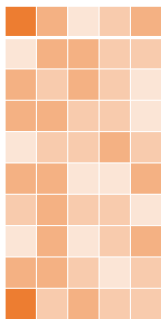
truth: $P \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$



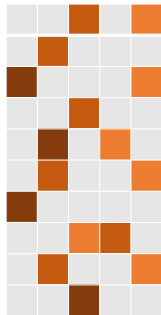
empirical estimate: \hat{P}

- 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?

l_∞ -based sample complexity

Theorem (Agarwal, Kakade, Yang '19)

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

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

with high prob., with sample complexity at most

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

ℓ_∞ -based sample complexity

Theorem (Agarwal, Kakade, Yang '19)

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

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

with high prob., with sample complexity at most

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

- matches minimax lower bound: $\tilde{\Omega}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$ when $\varepsilon \leq \frac{1}{\sqrt{1-\gamma}}$
(equivalently, when sample size exceeds $\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^2}$) Azar et al., 2013

ℓ_∞ -based sample complexity

Theorem (Agarwal, Kakade, Yang '19)

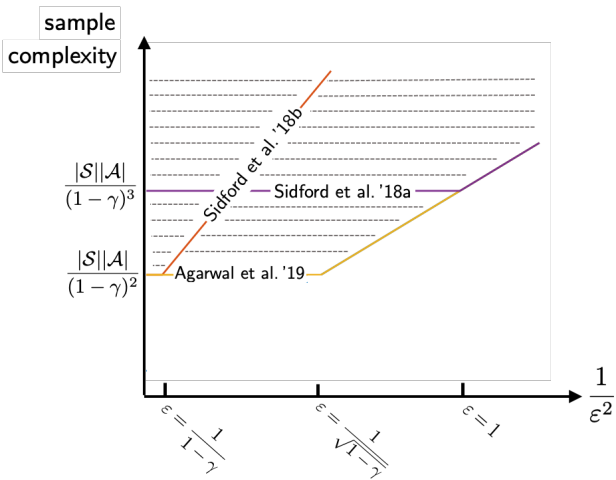
For any $0 < \varepsilon \leq \frac{1}{\sqrt{1-\gamma}}$, the optimal policy $\hat{\pi}^*$ of empirical MDP achieves

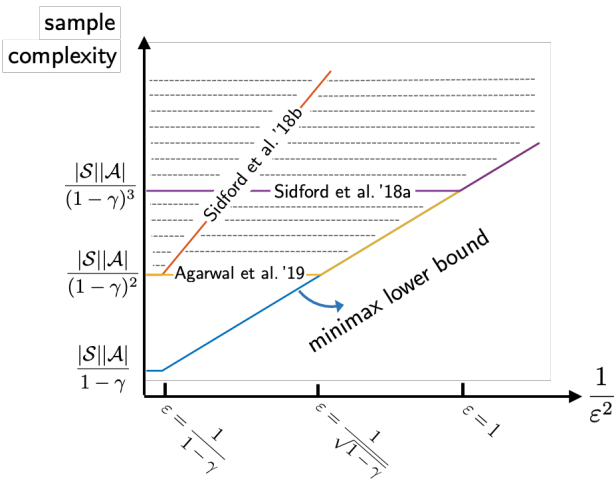
$$\|V^{\hat{\pi}^*} - V^*\|_\infty \leq \varepsilon$$

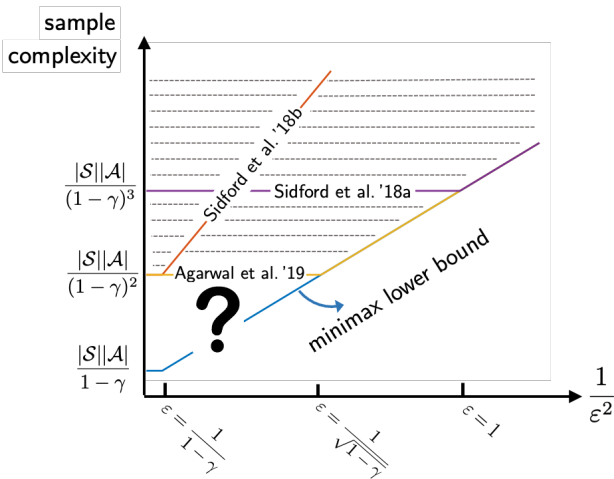
with high prob., with sample complexity at most

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

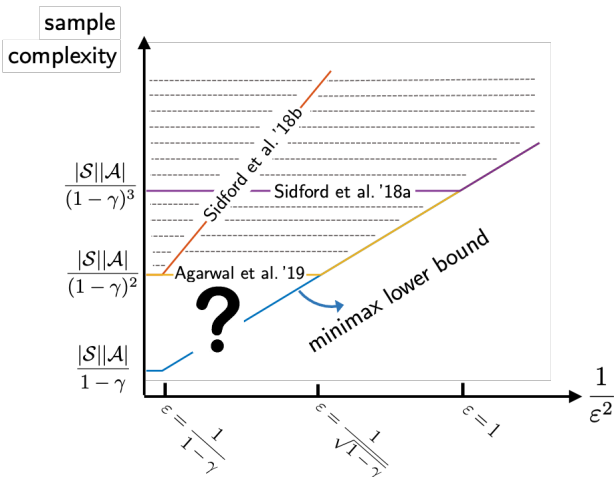
- matches minimax lower bound: $\tilde{\Omega}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$ when $\varepsilon \leq \frac{1}{\sqrt{1-\gamma}}$
(equivalently, when sample size exceeds $\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^2}$) Azar et al., 2013
- established upon leave-one-out analysis framework







Agarwal et al., 2019 still requires a burn-in sample size $\gtrsim \frac{|S||\mathcal{A}|}{(1-\gamma)^2}$

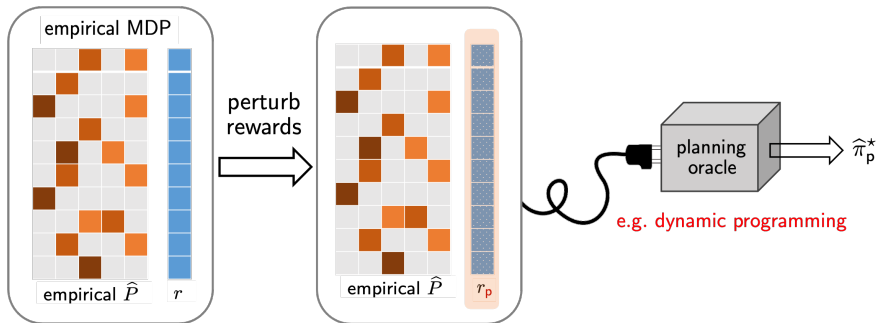


Agarwal et al., 2019 still requires a burn-in sample size $\gtrsim \frac{|S||\mathcal{A}|}{(1-\gamma)^2}$

Question: is it possible to break this sample size barrier?

Perturbed model-based approach (Li et al. '20)

—Li et al., 2020



Find policy based on the **empirical** MDP with **slightly perturbed** rewards

Optimal l_∞ -based sample complexity

Theorem (Li, Wei, Chi, Chen '20)

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

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

with high prob., with sample complexity at most

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

Optimal l_∞ -based sample complexity

Theorem (Li, Wei, Chi, Chen '20)

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

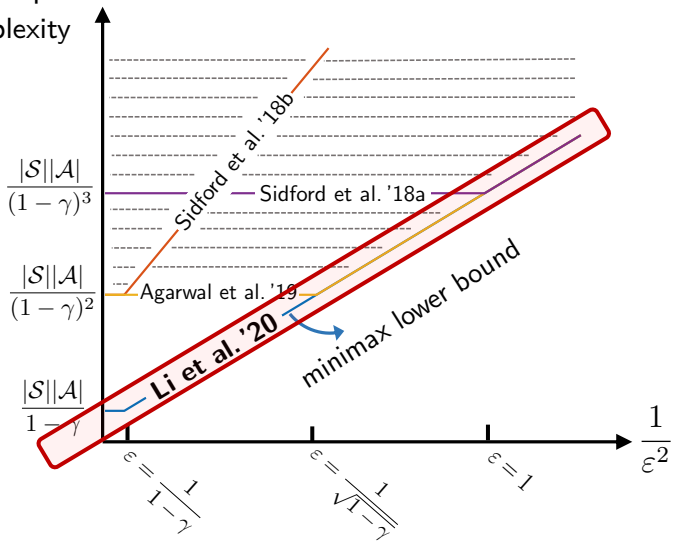
$$\|V^{\hat{\pi}_p^*} - V^*\|_\infty \leq \varepsilon$$

with high prob., with sample complexity at most

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

- matches minimax lower bound: $\tilde{\Omega}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$ [Azar et al., 2013](#)
- full ε -range: $\varepsilon \in (0, \frac{1}{1-\gamma}] \rightarrow$ no burn-in cost
- established upon more refined **leave-one-out analysis** and a perturbation argument

sample
complexity



Model-based RL (a “plug-in” approach)

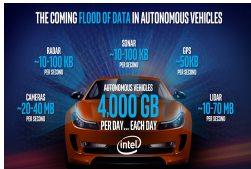
1. Sampling from a generative model (simulator)
2. Offline RL / batch RL
3. Robust RL

Offline RL / batch RL

- Collecting new data might be expensive or time-consuming
- But we have already stored tons of historical data



medical records



data of self-driving



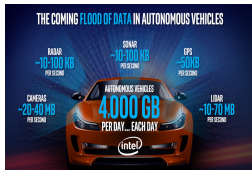
clicking times of ads

Offline RL / batch RL

- Collecting new data might be expensive or time-consuming
- But we have already stored tons of historical data



medical records



data of self-driving



clicking times of ads

Question: Can we design algorithms based solely on historical data?

Offline RL / batch RL

A historical dataset $\mathcal{D} = \{(s^{(i)}, a^{(i)}, s'^{(i)})\}$: N independent copies of

$$s \sim \rho^b, \quad a \sim \pi^b(\cdot | s), \quad s' \sim P(\cdot | s, a)$$

for some state distribution ρ^b and behavior policy π^b

Offline RL / batch RL

A historical dataset $\mathcal{D} = \{(s^{(i)}, a^{(i)}, s'^{(i)})\}$: N independent copies of

$$s \sim \rho^b, \quad a \sim \pi^b(\cdot | s), \quad s' \sim P(\cdot | s, a)$$

for some state distribution ρ^b and behavior policy π^b

Goal: given some test distribution ρ and accuracy level ε , find an ε -optimal policy $\hat{\pi}$ based on \mathcal{D} obeying

$$V^*(\rho) - V^{\hat{\pi}}(\rho) = \mathbb{E}_{s \sim \rho} [V^*(s)] - \mathbb{E}_{s \sim \rho} [V^{\hat{\pi}}(s)] \leq \varepsilon$$

— *in a sample-efficient manner*

Challenges of offline RL

- **Distribution shift:**

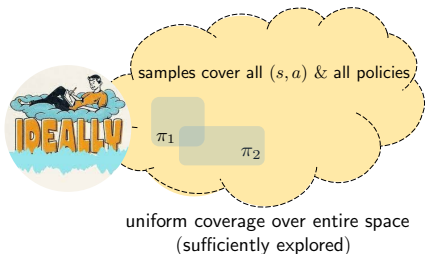
distribution(\mathcal{D}) \neq target distribution under π^*

Challenges of offline RL

- **Distribution shift:**

distribution(\mathcal{D}) \neq target distribution under π^*

- **Partial coverage of state-action space:**

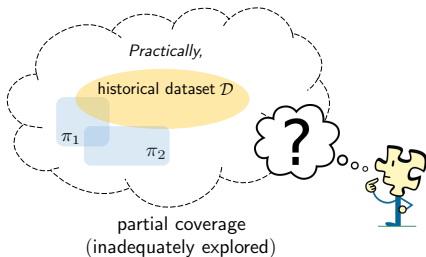
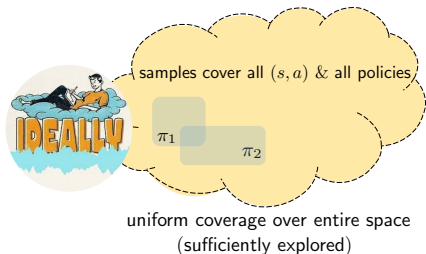


Challenges of offline RL

- **Distribution shift:**

distribution(\mathcal{D}) \neq target distribution under π^*

- **Partial coverage of state-action space:**



How to quantify quality of historical dataset \mathcal{D} (induced by π^b)?

How to quantify quality of historical dataset \mathcal{D} (induced by π^b)?

Single-policy concentrability coefficient

$$C^* := \max_{s,a} \frac{d^{\pi^*}(s,a)}{d^{\pi^b}(s,a)}$$

where $d^\pi(s,a) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}((s^t, a^t) = (s,a) | \pi)$

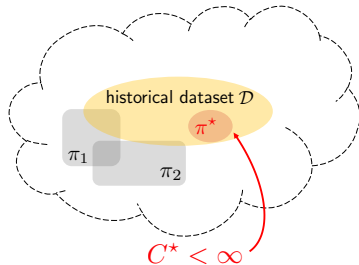
How to quantify quality of historical dataset \mathcal{D} (induced by π^b)?

Single-policy concentrability coefficient

$$C^* := \max_{s,a} \frac{d^{\pi^*}(s,a)}{d^{\pi^b}(s,a)} = \left\| \frac{\text{occupancy density of } \pi^*}{\text{occupancy density of } \pi^b} \right\|_{\infty} \geq 1$$

where $d^{\pi}(s,a) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}((s^t, a^t) = (s,a) | \pi)$

- captures distributional shift
- allows for partial coverage



Key idea: pessimism in the face of uncertainty

— *Jin et al. '20, Rashidinejad et al. '21, Xie et al. '21*



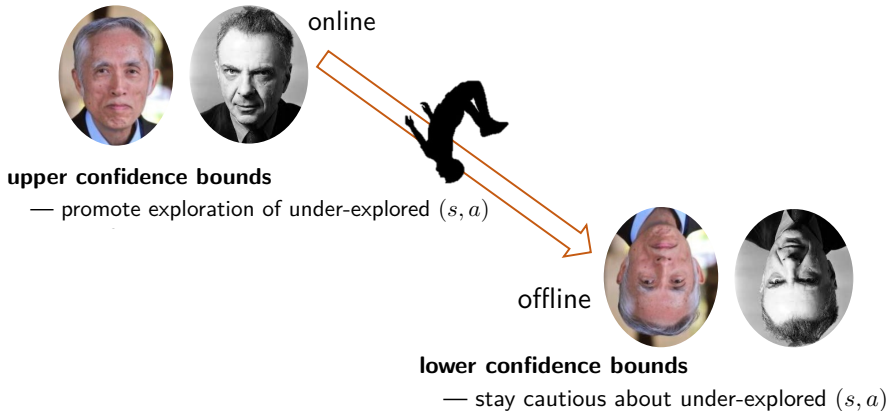
online

upper confidence bounds

— promote exploration of under-explored (s, a)

Key idea: pessimism in the face of uncertainty

— Jin et al. '20, Rashidinejad et al. '21, Xie et al. '21



Key idea: pessimism in the face of uncertainty

— Jin et al. '20, Rashidinejad et al. '21, Xie et al. '21

A model-based offline algorithm: VI-LCB

1. build empirical model \hat{P}
2. **(value iteration)** for $t \leq \tau_{\max}$:

$$\hat{Q}_t(s, a) \leftarrow \left[r(s, a) + \gamma \langle \hat{P}(\cdot | s, a), \hat{V}_{t-1} \rangle \right]_+$$

for all (s, a) , where $\hat{V}_t(s) = \max_a \hat{Q}_t(s, a)$

Key idea: pessimism in the face of uncertainty

— Jin et al. '20, Rashidinejad et al. '21, Xie et al. '21

A model-based offline algorithm: VI-LCB

1. build empirical model \hat{P}
2. **(pessimistic value iteration)** for $t \leq \tau_{\max}$:

$$\hat{Q}_t(s, a) \leftarrow \left[r(s, a) + \gamma \langle \hat{P}(\cdot | s, a), \hat{V}_{t-1} \rangle - \underbrace{b(s, a; \hat{V}_{t-1})}_{\text{penalize poorly visited } (s, a)} \right]_+$$

for all (s, a) , where $\hat{V}_t(s) = \max_a \hat{Q}_t(s, a)$

Key idea: pessimism in the face of uncertainty

— Jin et al. '20, Rashidinejad et al. '21, Xie et al. '21

A model-based offline algorithm: VI-LCB

1. build empirical model \hat{P}
2. **(pessimistic value iteration)** for $t \leq \tau_{\max}$:

$$\hat{Q}_t(s, a) \leftarrow \left[r(s, a) + \gamma \langle \hat{P}(\cdot | s, a), \hat{V}_{t-1} \rangle - \underbrace{b(s, a; \hat{V}_{t-1})}_{\text{penalize poorly visited } (s, a)} \right]_+$$

compared w/ prior works

- no need of variance reduction
- variance-aware penalty

Minimax optimality of model-based offline RL

Theorem (Li, Shi, Chen, Chi, Wei '22)

For any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the policy $\hat{\pi}$ returned by VI-LCB achieves

$$V^*(\rho) - V^{\hat{\pi}}(\rho) \leq \varepsilon$$

with high prob., with sample complexity at most

$$\tilde{O}\left(\frac{SC^*}{(1-\gamma)^3\varepsilon^2}\right)$$

Minimax optimality of model-based offline RL

Theorem (Li, Shi, Chen, Chi, Wei '22)

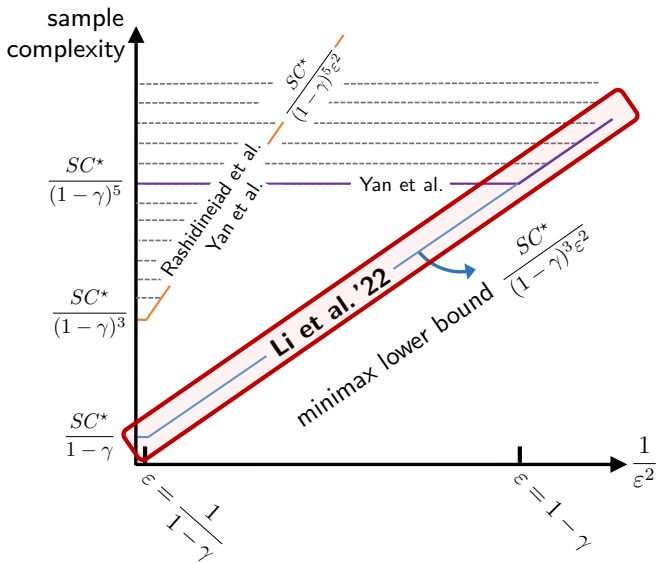
For any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the policy $\hat{\pi}$ returned by VI-LCB achieves

$$V^*(\rho) - V^{\hat{\pi}}(\rho) \leq \varepsilon$$

with high prob., with sample complexity at most

$$\tilde{O}\left(\frac{SC^*}{(1-\gamma)^3\varepsilon^2}\right)$$

- matches minimax lower bound: $\tilde{\Omega}\left(\frac{SC^*}{(1-\gamma)^3\varepsilon^2}\right)$ [Rashidinejad et al, 2021](#)
- depends on distribution shift (as reflected by C^*)
- full ε -range (no burn-in cost)

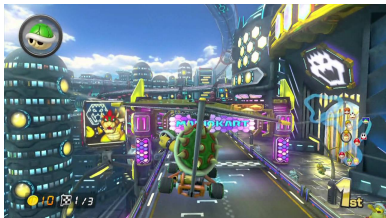


Model-based RL (a “plug-in” approach)

1. Sampling from a generative model (simulator)
2. Offline RL / batch RL
3. Robust RL

Safety and robustness in RL

(Zhou et al., 2021; Panaganti and Kalathil, 2022; Yang et al., 2022;)



Training environment

≠



Test environment

Safety and robustness in RL

(Zhou et al., 2021; Panaganti and Kalathil, 2022; Yang et al., 2022;)



Training environment

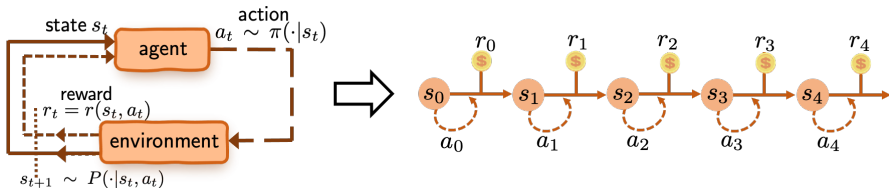
≠



Test environment

Sim2Real Gap: Can we learn optimal policies that are robust to model perturbations?

Distributionally robust MDP



Uncertainty set of the nominal transition kernel P^o :

$$\mathcal{U}^\sigma(P^o) = \{P : \rho(P, P^o) \leq \sigma\}$$

Robust value/Q function of policy π :

$$\forall s \in \mathcal{S} : \quad V^{\pi, \sigma}(s) := \inf_{P \in \mathcal{U}^\sigma(P^o)} \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s \right]$$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad Q^{\pi, \sigma}(s, a) := \inf_{P \in \mathcal{U}^\sigma(P^o)} \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a \right]$$

The optimal robust policy π^* maximizes $V^{\pi, \sigma}(\rho)$

Robust Bellman's optimality equation

(Iyengar. '05, Nilim and El Ghaoui. '05)

Robust Bellman's optimality equation: the optimal robust policy π^* and optimal robust value $V^{*,\sigma} := V^{\pi^*,\sigma}$ satisfy

$$Q^{*,\sigma}(s, a) = r(s, a) + \gamma \inf_{P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^o)} \langle P_{s,a}, V^{*,\sigma} \rangle,$$

$$V^{*,\sigma}(s) = \max_a Q^{*,\sigma}(s, a)$$

Robust Bellman's optimality equation

(Iyengar. '05, Nilim and El Ghaoui. '05)

Robust Bellman's optimality equation: the optimal robust policy π^* and optimal robust value $V^{*,\sigma} := V^{\pi^*,\sigma}$ satisfy

$$Q^{*,\sigma}(s, a) = r(s, a) + \gamma \inf_{P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^o)} \langle P_{s,a}, V^{*,\sigma} \rangle,$$

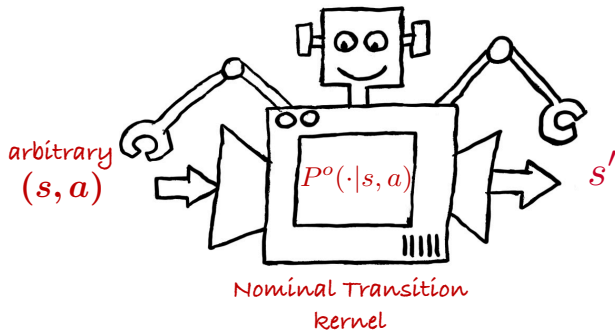
$$V^{*,\sigma}(s) = \max_a Q^{*,\sigma}(s, a)$$

Robust value iteration:

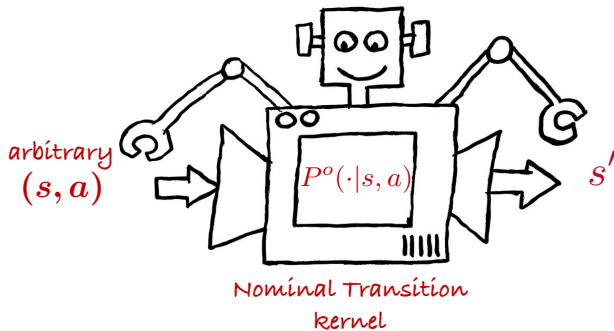
$$Q(s, a) \leftarrow r(s, a) + \gamma \inf_{P_{s,a} \in \mathcal{U}^\sigma(P_{s,a}^o)} \langle P_{s,a}, V \rangle,$$

where $V(s) = \max_a Q(s, a)$.

Learning distributionally robust MDPs



Learning distributionally robust MDPs

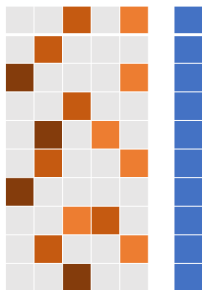


Goal of robust RL: given $\mathcal{D} := \{(s_i, a_i, s'_i)\}_{i=1}^N$ from the *nominal* environment P^0 , find an ε -optimal robust policy $\hat{\pi}$ obeying

$$V^{*,\sigma}(\rho) - V^{\hat{\pi},\sigma}(\rho) \leq \varepsilon$$

— in a sample-efficient manner

A curious question



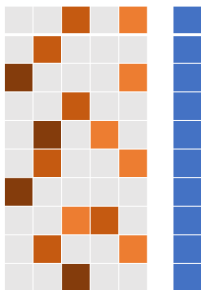
empirical MDP

Learn the optimal policy of the nominal MDP?

Learn the **robust** policy around the nominal MDP?



A curious question



empirical MDP

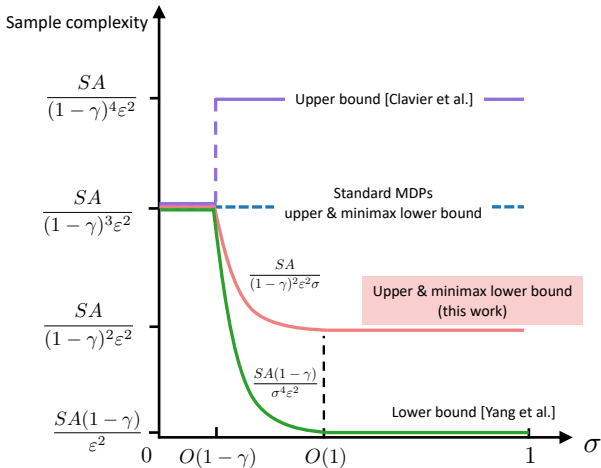
Learn the optimal policy of the nominal MDP?

Learn the **robust** policy around the nominal MDP?

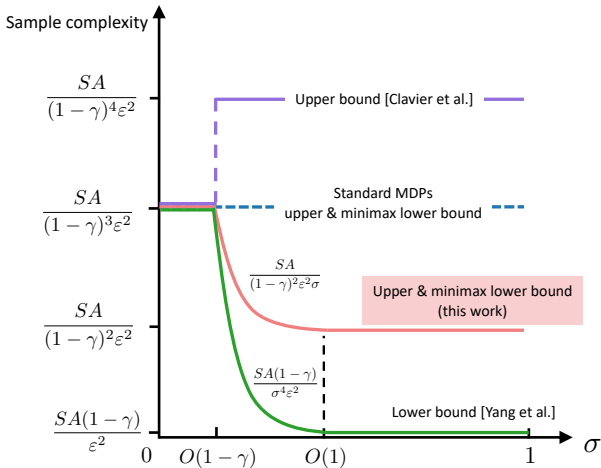


Robustness-statistical trade-off? Is there a statistical premium that one needs to pay in quest of additional robustness?

When the uncertainty set is TV

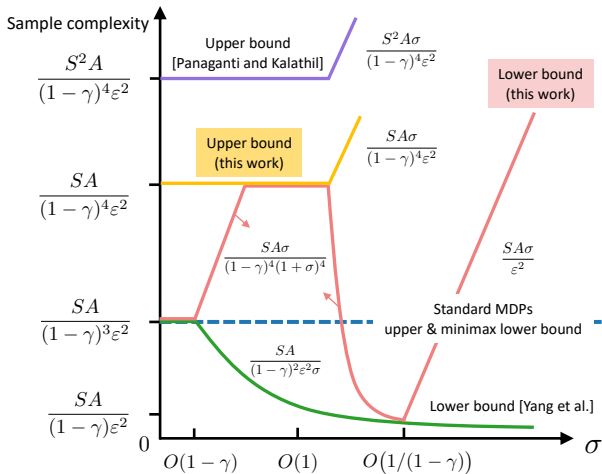


When the uncertainty set is TV

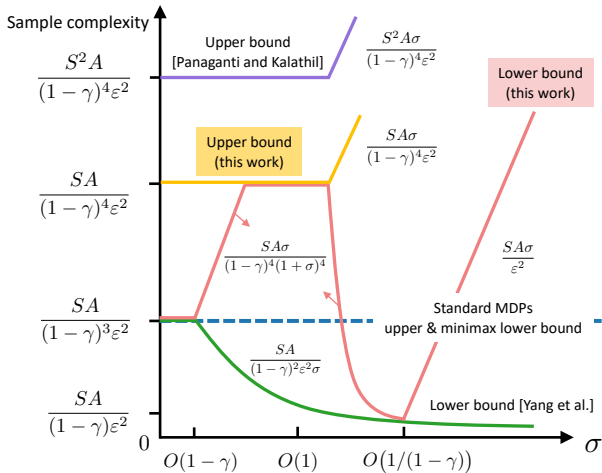


RMDPs are **easier** to learn than standard MDPs.

When the uncertainty set is χ^2 divergence



When the uncertainty set is χ^2 divergence



RMDPs can be **harder** to learn than standard MDPs.

Summary of this part

Model-based RL (a “plug-in” approach)

- Sampling from a generative model (simulator)
- Offline RL / batch RL
- Robust RL

Papers:

“Breaking the sample size barrier in model-based reinforcement learning with a generative model,” G Li, Y Wei, Y Chi, Y Chen, *NeurIPS’20, Operators Research’23*

“Settling the sample complexity of model-based offline reinforcement learning,” G Li, L Shi, Y Chen, Y Chi, Y Wei, 2022

“The curious price of distributional robustness in reinforcement learning with a generative model,” L Shi, G Li, Y Wei, Y Chen, M Geist, Y Chi, 2023