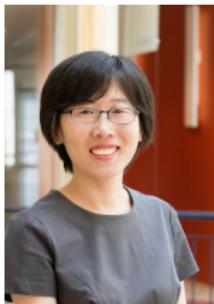


Recent Advances of Statistical Reinforcement Learning

Part 1



Yuejie Chi
Carnegie Mellon University



Sattar Vakili
MediaTek Research



Gergely Neu
Universitat Pompeu Fabra

Tutorial, UAI 2024

Acknowledgement

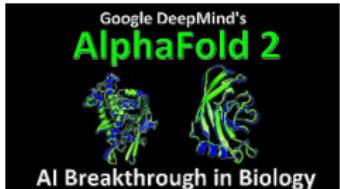


Yuxin Chen
UPenn



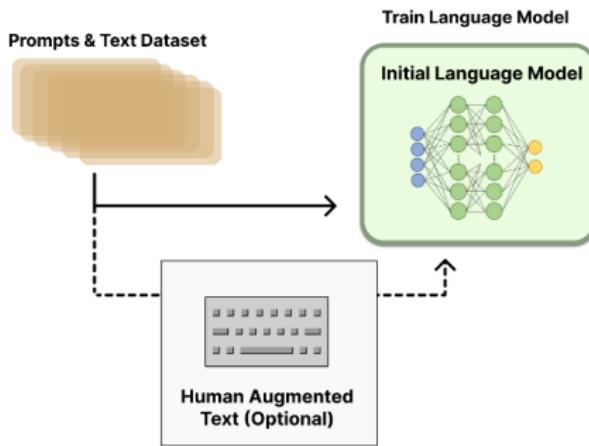
Yuting Wei
UPenn

Recent successes in reinforcement learning (RL)



RL holds great promise in the era of AI

One more recent success: RLHF



You

What is RLHF?



ChatGPT

RLHF stands for **Reinforcement Learning from Human Feedback**. It's a technique used in machine learning and artificial intelligence where a model learns to perform tasks or make decisions based on feedback from human trainers, rather than solely relying on pre-existing data sets or explicit programming. This approach allows the

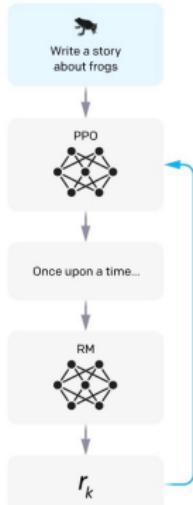


ChatGPT



Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.



The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

Data efficiency

Data collection might be expensive, time-consuming, or high-stakes



clinical trials



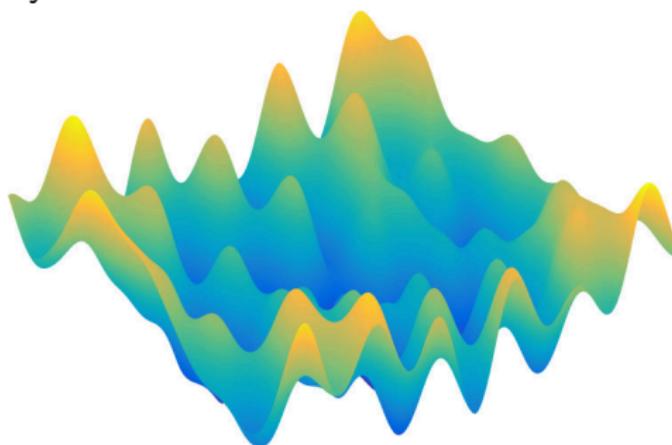
self-driving cars

Calls for design of sample-efficient RL algorithms!

Computational efficiency

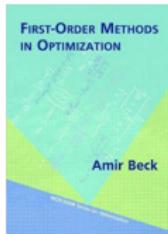
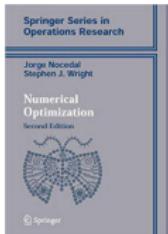
Running RL algorithms might take a long time ...

- enormous state-action space
- nonconvexity

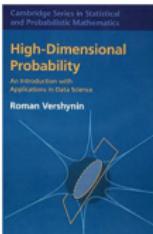
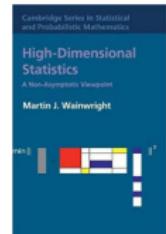


Calls for computationally efficient RL algorithms!

This tutorial



(large-scale) optimization



(high-dimensional) statistics

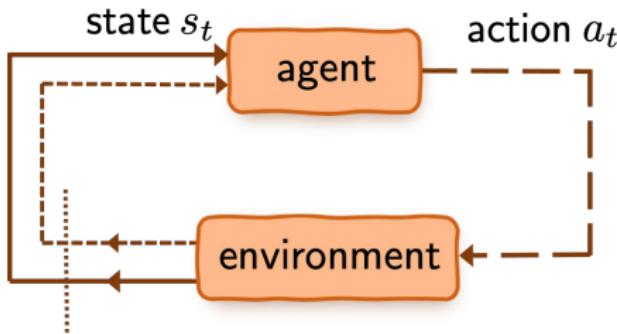
Part 1. Basics, statistical RL in the tabular setting

Part 2. Beyond the tabular setting

Part 1

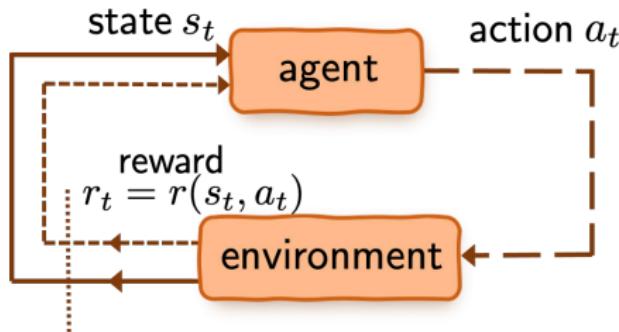
1. Basics: Markov decision processes
2. RL w/ a generative model (simulator)
3. Online RL
4. Offline RL

Markov decision process (MDP)



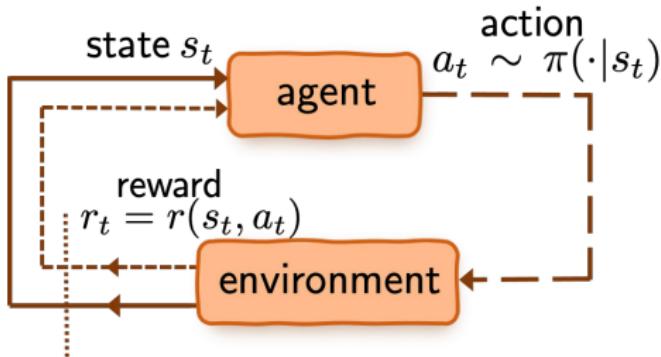
- $\mathcal{S} = \{1, \dots, S\}$: state space (containing S states)
- $\mathcal{A} = \{1, \dots, A\}$: action space (containing A actions)

Markov decision process (MDP)



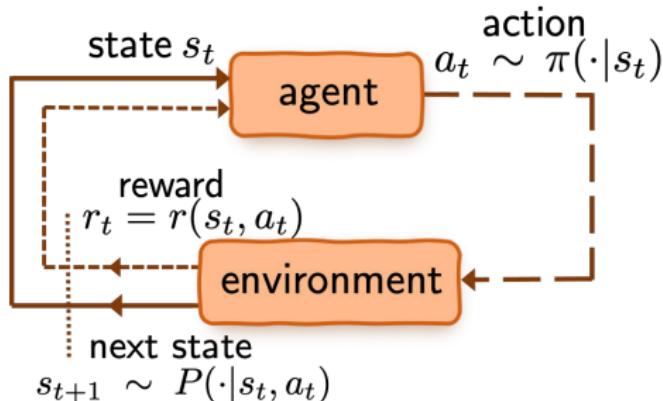
- $\mathcal{S} = \{1, \dots, S\}$: state space (containing S states)
- $\mathcal{A} = \{1, \dots, A\}$: action space (containing A actions)
- $r(s, a) \in [0, 1]$: immediate reward

Discounted infinite-horizon MDPs



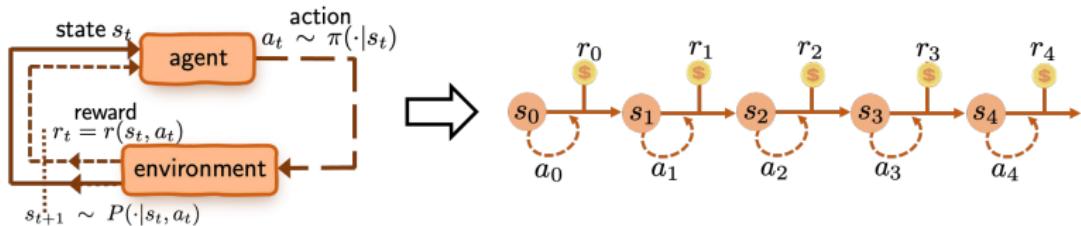
- $\mathcal{S} = \{1, \dots, S\}$: state space (containing S states)
- $\mathcal{A} = \{1, \dots, A\}$: action space (containing A actions)
- $r(s, a) \in [0, 1]$: immediate reward
- $\pi(\cdot | s)$: policy (or action selection rule)

Discounted infinite-horizon MDPs



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

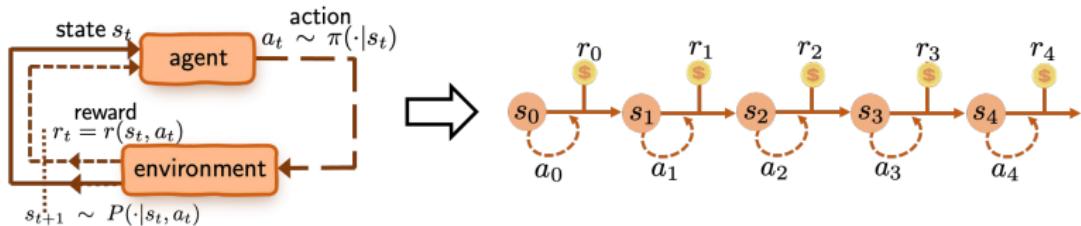
Value function



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

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

Value function

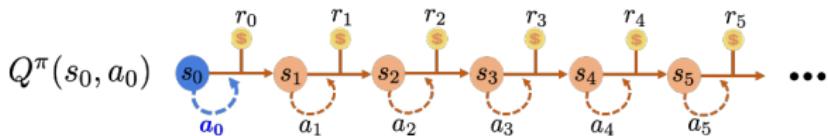


Value of policy π : cumulative **discounted** reward

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

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

Q-function (action-value function)

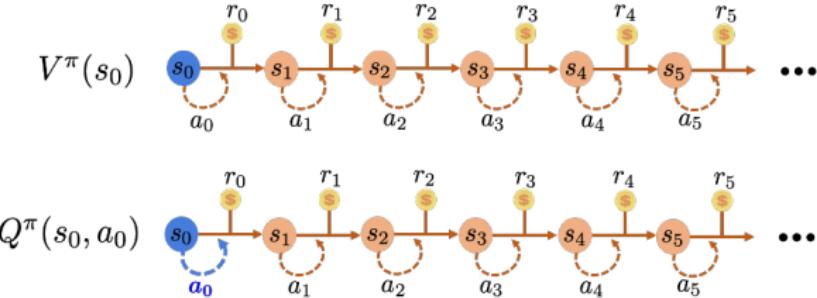


Q-function of policy π :

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

- $(\textcolor{red}{a_0}, s_1, a_1, s_2, a_2, \dots)$: induced by policy π

Q-function (action-value function)

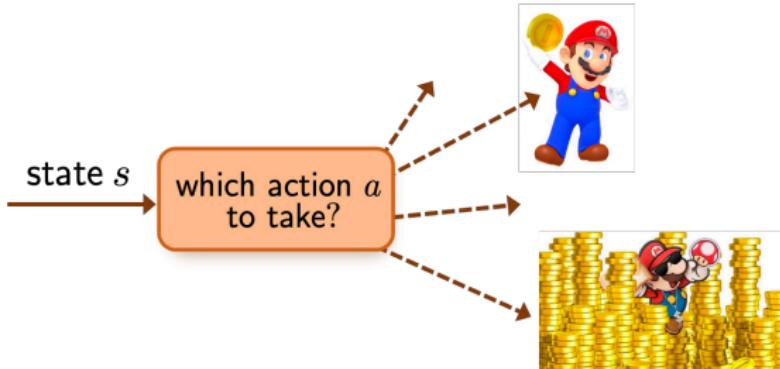


Q-function of policy π :

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

- $(\textcolor{red}{a_0}, s_1, a_1, s_2, a_2, \dots)$: induced by policy π

Optimal policy and optimal value



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

Optimal policy and optimal value



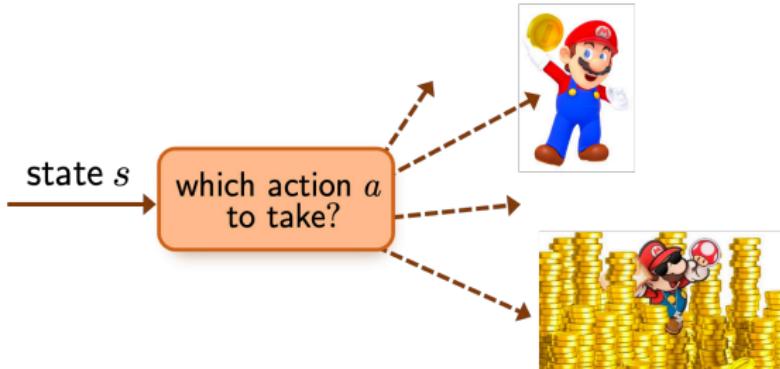
- **optimal policy** π^* : maximizing value function $\max_{\pi} V^{\pi}$

Theorem 1 (Puterman'94)

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

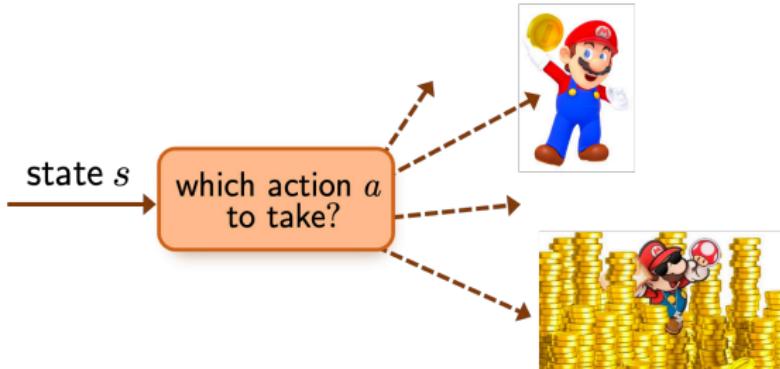
$$V^{\pi^*}(s) \geq V^{\pi}(s), \quad \forall s, \text{ and } \pi.$$

Optimal policy and optimal value



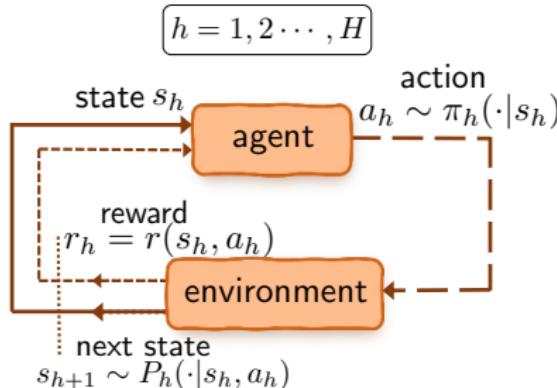
- **optimal policy** π^* : maximizing value function $\max_{\pi} V^{\pi}$
- **optimal value / Q function**: $V^* := V^{\pi^*}$, $Q^* := Q^{\pi^*}$

Optimal policy and optimal value



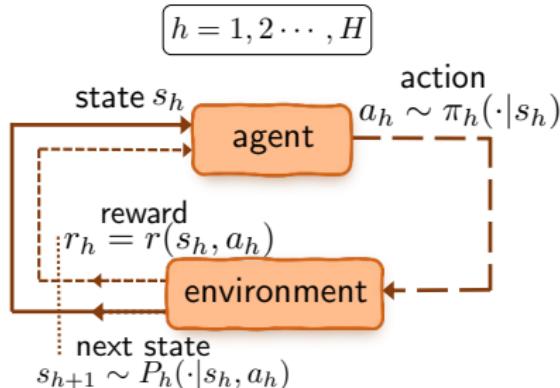
- **optimal policy** π^* : maximizing value function $\max_{\pi} V^{\pi}$
- **optimal value / Q function**: $V^* := V^{\pi^*}$, $Q^* := Q^{\pi^*}$
- A question to keep in mind: *how to find optimal π^* ?*

Finite-horizon MDPs (nonstationary)



- 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 (nonstationary)

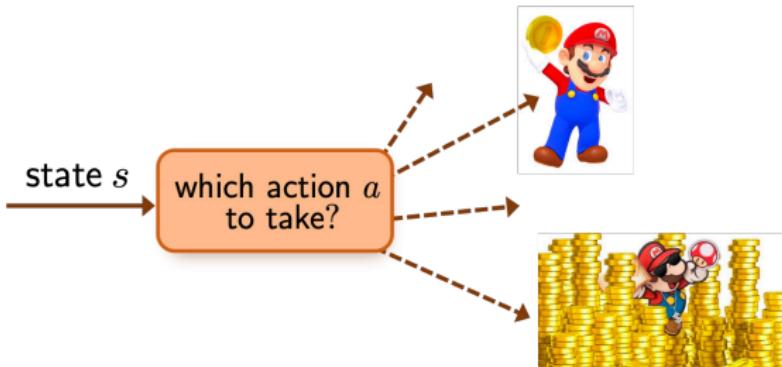


value function: $V_h^\pi(s) := \mathbb{E} \left[\sum_{t=h}^H r_h(s_h, a_h) \mid s_h = s \right]$

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



Optimal policy and optimal value



- **optimal policy** π^* : maximizing value function at all steps
- **optimal value / Q function**: $V_h^* := V_h^{\pi^*}$, $Q_h^* := Q_h^{\pi^*}$, $\forall h$
- **Question:** *how to find optimal π^* ?*

*Basic dynamic programming algorithms
when MDP specification is known*

A simpler problem: **policy evaluation**

— given MDP \mathcal{M} and policy π , how to compute V^π , Q^π ?

A simpler problem: **policy evaluation**

— given MDP \mathcal{M} and policy π , how to compute V^π , Q^π ?

A simpler problem: **policy evaluation**

— given MDP \mathcal{M} and policy π , how to compute V^π , Q^π ?

solution: 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

A simpler problem: **policy evaluation**

— given MDP \mathcal{M} and policy π , how to compute V^π , Q^π ?

solution: 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
- P^π : state-action transition matrix induced by π :

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



Richard Bellman

Back to main question: how to find optimal policy π^* ?

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

Back to main question: how to find optimal policy π^* ?

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

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

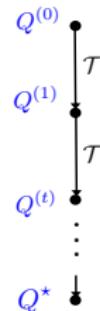
$$\mathcal{T}(Q^*) = Q^*$$

Two dynamic programming algorithms

Value iteration (VI)

For $t = 0, 1, \dots$

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

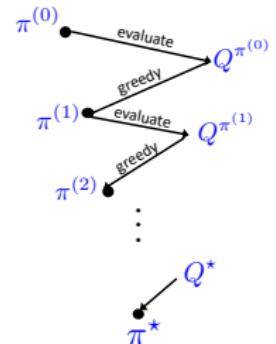


Policy iteration (PI)

For $t = 0, 1, \dots$

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

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



Iteration complexity

Theorem 2 (Linear convergence of policy/value iteration)

$$\|Q^{(t)} - Q^*\|_\infty \leq \gamma^t \|Q^{(0)} - Q^*\|_\infty$$

Iteration complexity

Theorem 2 (Linear convergence of policy/value iteration)

$$\|Q^{(t)} - Q^*\|_\infty \leq \gamma^t \|Q^{(0)} - Q^*\|_\infty$$

Implications: to achieve $\|Q^{(t)} - Q^*\|_\infty \leq \varepsilon$, it takes no more than

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

Iteration complexity

Theorem 2 (Linear convergence of policy/value iteration)

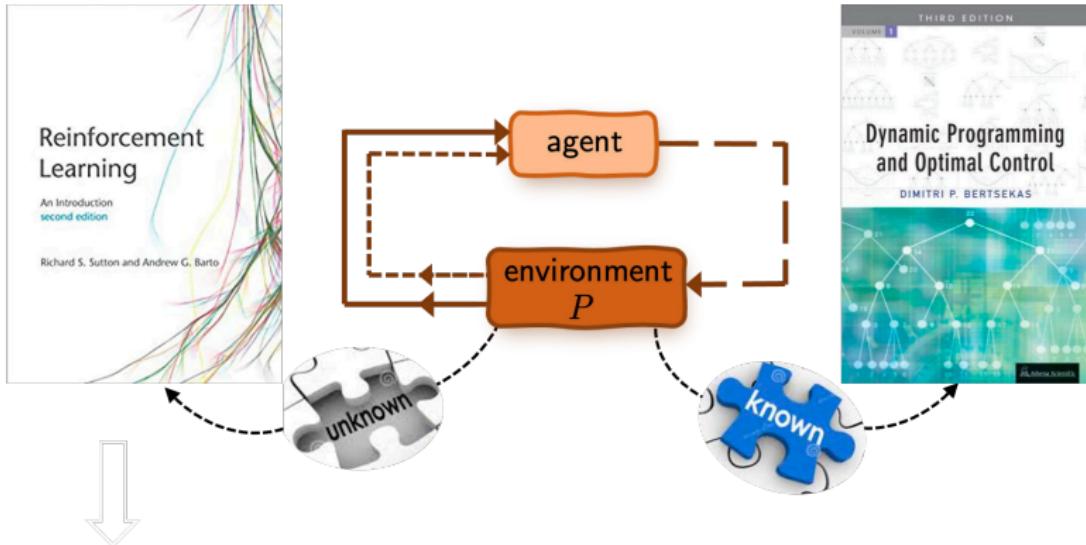
$$\|Q^{(t)} - Q^*\|_\infty \leq \gamma^t \|Q^{(0)} - Q^*\|_\infty$$

Implications: to achieve $\|Q^{(t)} - Q^*\|_\infty \leq \varepsilon$, it takes no more than

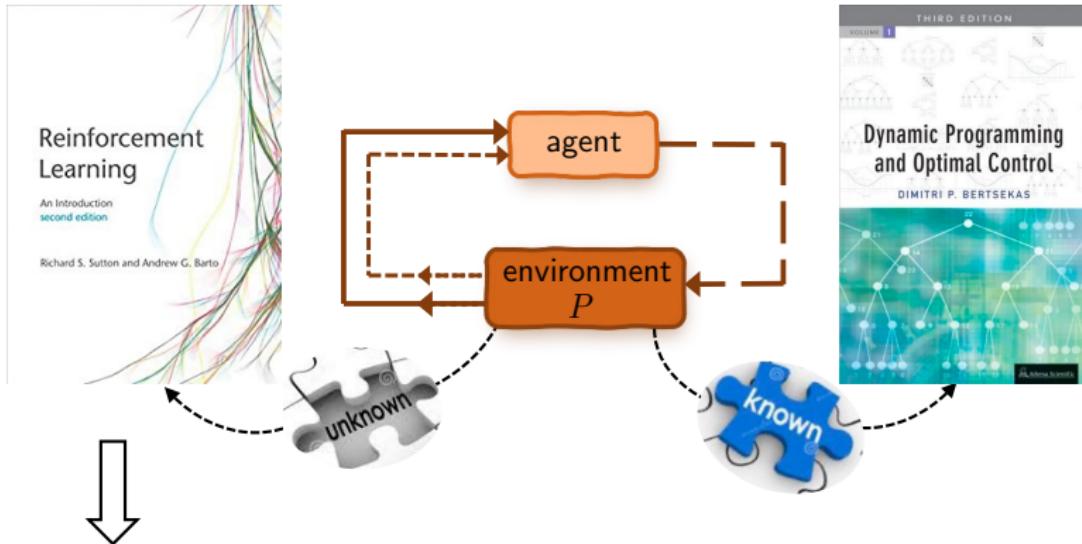
$$\frac{1}{1-\gamma} \log \left(\frac{\|Q^{(0)} - Q^*\|_\infty}{\varepsilon} \right) \text{ iterations}$$

Linear convergence at a **dimension-free** rate!

When the model is unknown ...

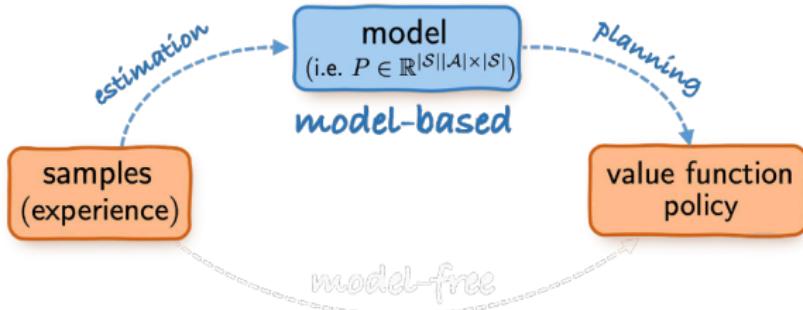


When the model is unknown ...



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

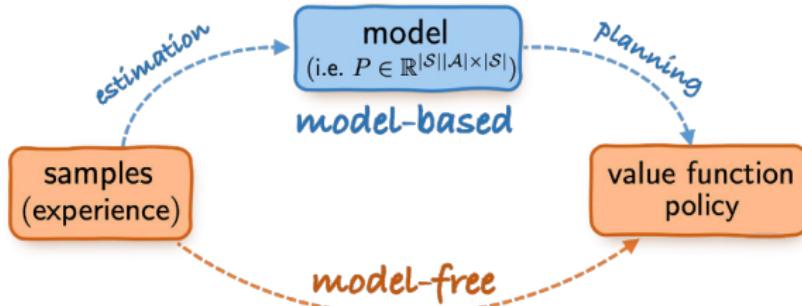
Two approaches



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

- learning w/o modeling & estimating environment explicitly
- memory-efficient, online, ...

Sampling mechanisms

1. RL w/ a generative model (a.k.a. simulator)
 - o can query arbitrary state-action pairs to draw samples

Sampling mechanisms

1. RL w/ a generative model (a.k.a. simulator)
 - can query arbitrary state-action pairs to draw samples
2. online RL
 - execute MDP in real time to obtain sample trajectories

Sampling mechanisms

1. RL w/ a generative model (a.k.a. simulator)
 - can query arbitrary state-action pairs to draw samples
2. online RL
 - execute MDP in real time to obtain sample trajectories
3. offline RL
 - use pre-collected historical data

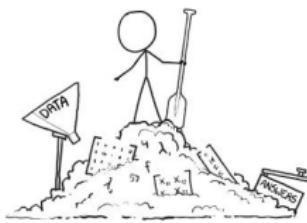
Sampling mechanisms

1. RL w/ a generative model (a.k.a. simulator)
 - can query arbitrary state-action pairs to draw samples
2. online RL
 - execute MDP in real time to obtain sample trajectories
3. offline RL
 - use pre-collected historical data

Question: how many samples are sufficient to learn an ε -optimal policy?

$$\hat{V^\pi} \geq V^* - \varepsilon$$

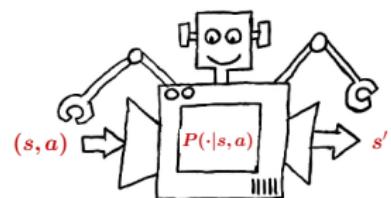
Exploration vs exploitation



offline RL

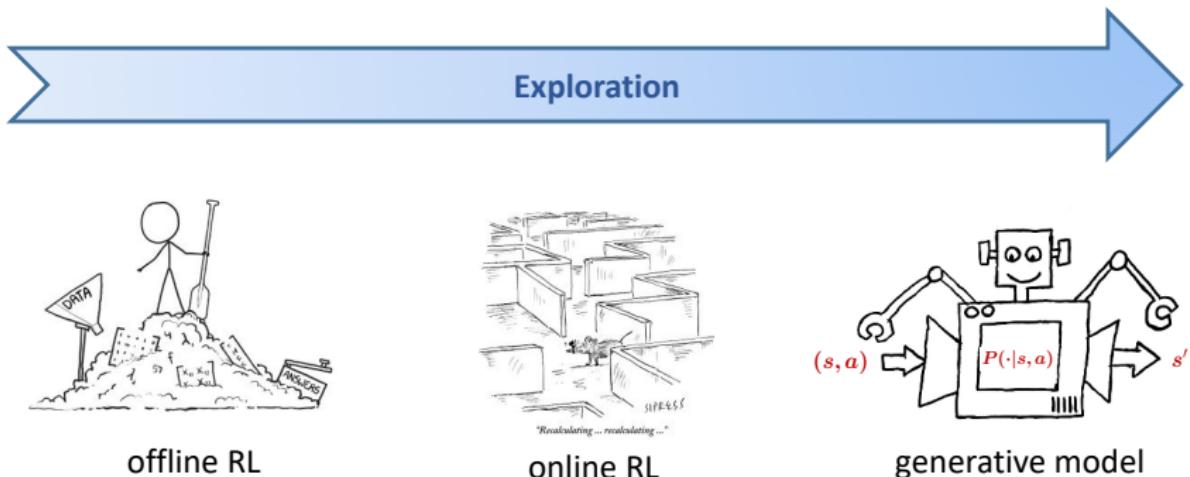


online RL



generative model

Exploration vs exploitation



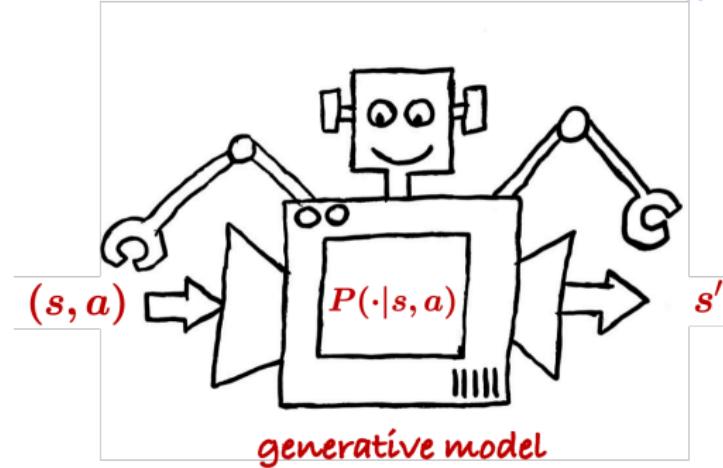
Varying levels of trade-offs between exploration and exploitation.

Part 1

1. Basics: Markov decision processes
2. RL w/ a generative model (simulator)
3. Online RL
4. Offline RL

A generative model / simulator

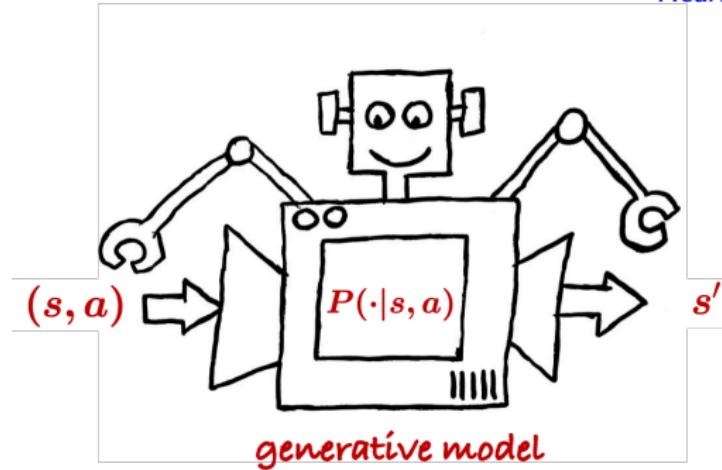
— Kearns and Singh, 1999



- **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 $SA \times N$)

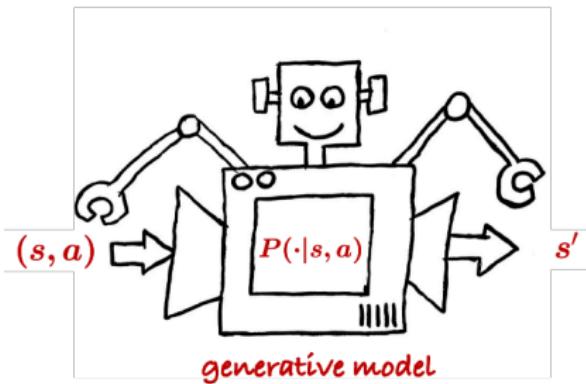
ℓ_∞ -sample complexity: how many samples are required to
learn an ε -optimal policy ?

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

An incomplete list of works

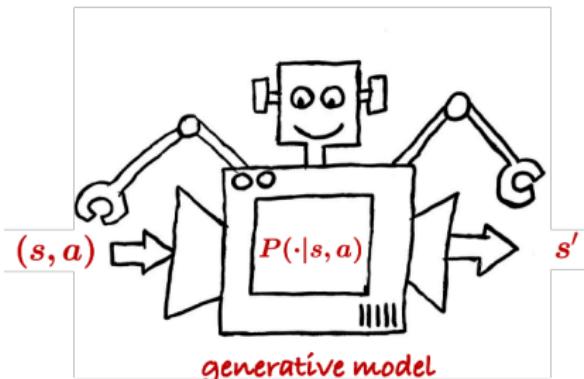
- Kearns and Singh, 1999
- Kakade, 2003
- Kearns et al., 2002
- Azar et al., 2013
- Sidford et al., 2018a, 2018b
- Wang, 2019
- Agarwal et al., 2019
- Wainwright, 2019a, 2019b
- Pananjady and Wainwright, 2019
- Yang and Wang, 2019
- Khamaru et al., 2020
- Mou et al., 2020
- Cui and Yang, 2021
- ...

Model estimation



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

Model estimation



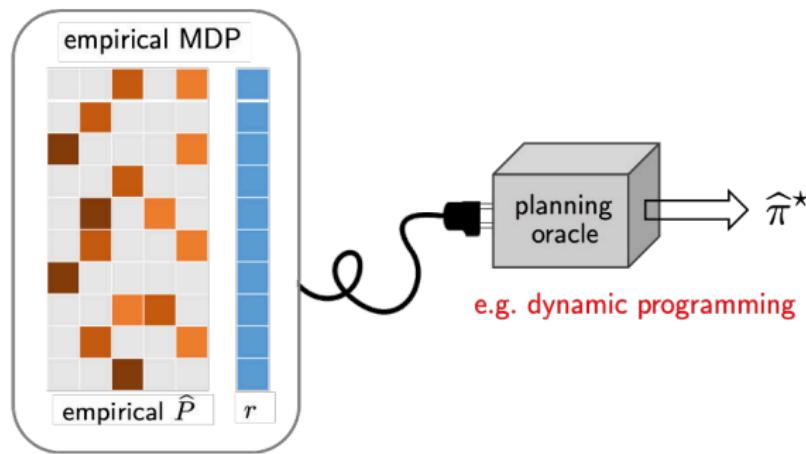
Sampling: for each (s, a) ,
collect N ind. samples
 $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$

Empirical estimates:

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

Empirical MDP + planning

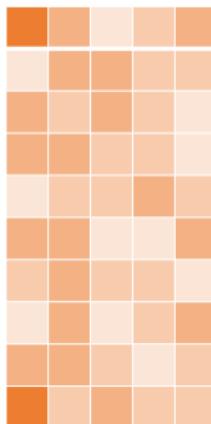
— Azar et al., 2013, Agarwal et al., 2019



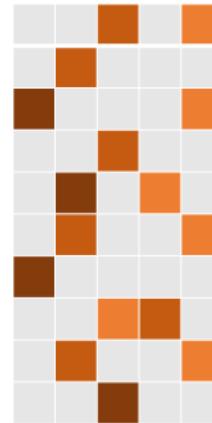
Find policy based on the empirical MDP (*empirical maximizer*)
using, e.g., policy iteration

$$(\hat{P}, r)$$

Challenges in the sample-starved regime



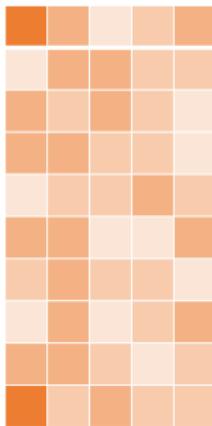
truth: $P \in \mathbb{R}^{SA \times S}$



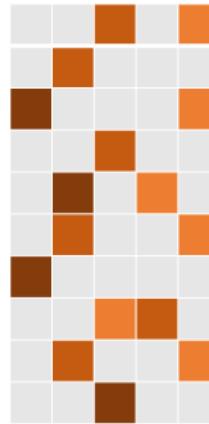
empirical estimate:
 \hat{P}

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

Challenges in the sample-starved regime



truth: $P \in \mathbb{R}^{SA \times S}$



empirical estimate:
 \hat{P}

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

ℓ_∞ -based sample complexity

Theorem 3 (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

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

ℓ_∞ -based sample complexity

Theorem 3 (Agarwal, Kakade, Yang '19)

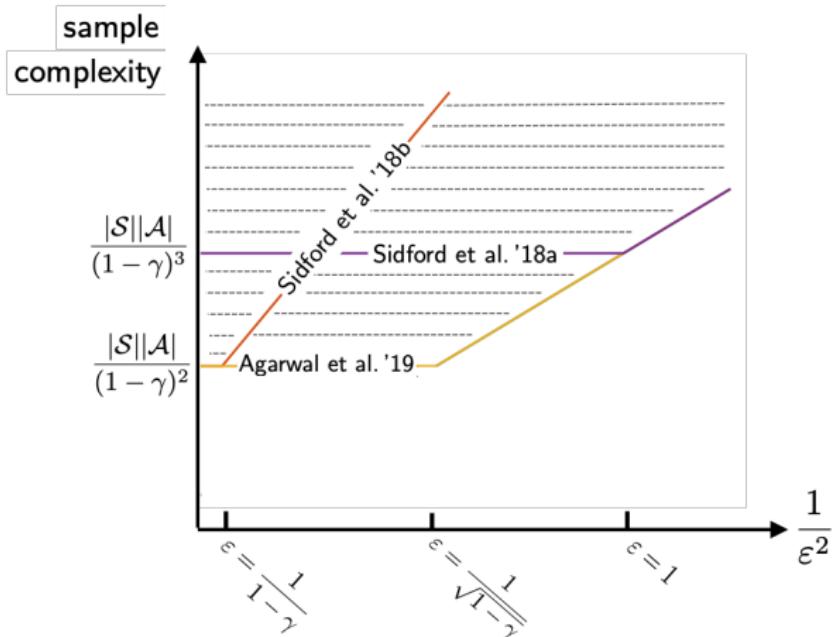
For any $0 < \varepsilon \leq \frac{1}{\sqrt{1-\gamma}}$, the optimal policy $\hat{\pi}^*$ of empirical MDP achieves

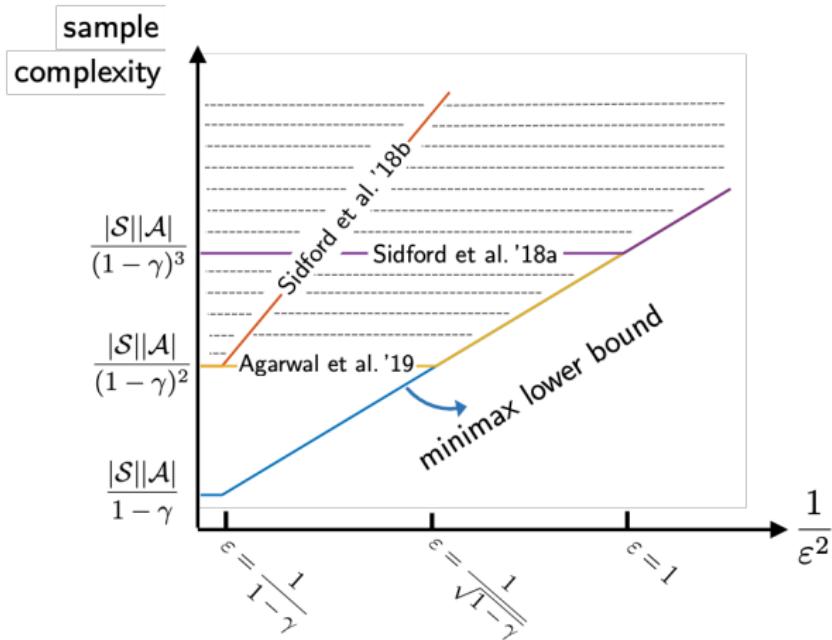
$$\|V^{\hat{\pi}^*} - V^*\|_\infty \leq \varepsilon$$

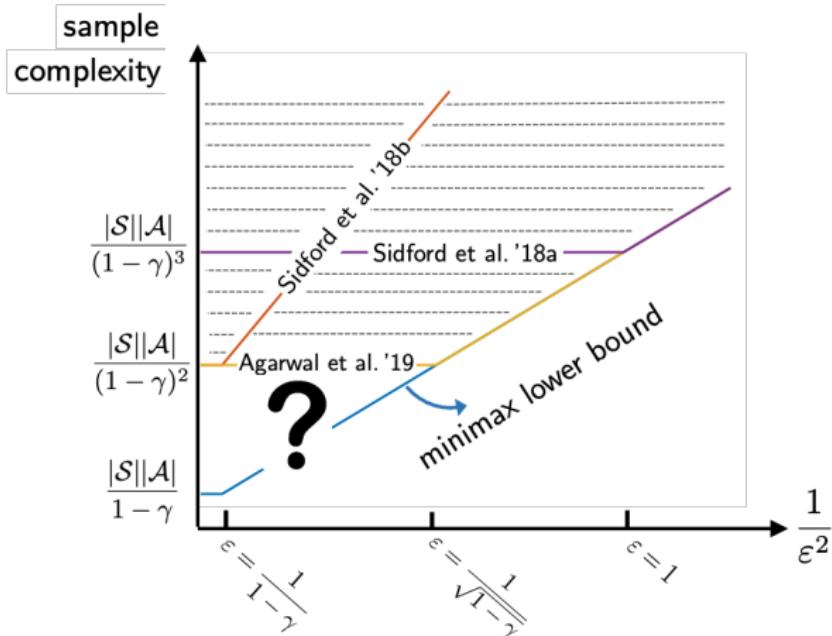
with high prob., with sample complexity at most

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

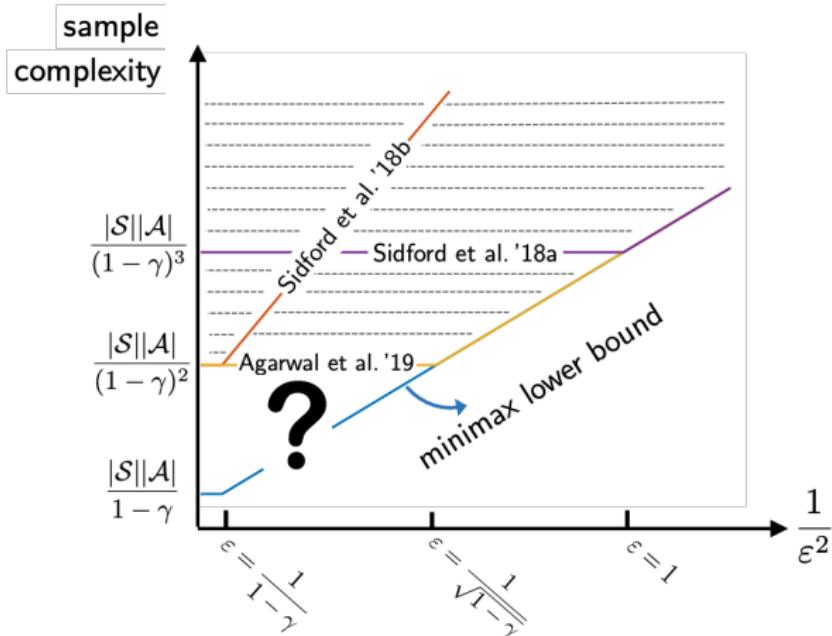
- matches minimax lower bound: $\tilde{\Omega}\left(\frac{SA}{(1-\gamma)^3\varepsilon^2}\right)$ when $\varepsilon \leq \frac{1}{\sqrt{1-\gamma}}$
(equivalently, when sample size exceeds $\frac{SA}{(1-\gamma)^2}$) Azar et al., 2013







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

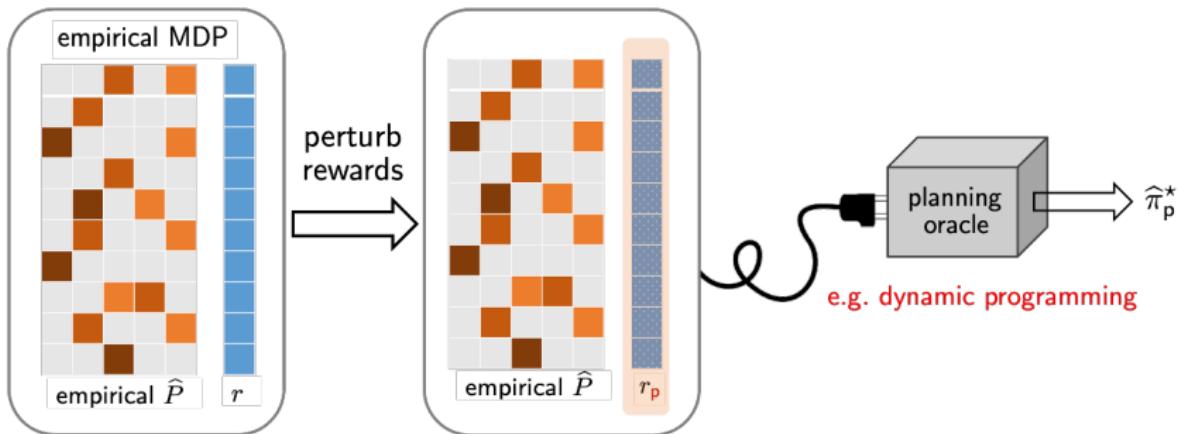


Agarwal et al., 2019 still requires a **burn-in sample size** $\gtrsim \frac{SA}{(1-\gamma)^2}$

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

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

— Li, Wei, Chi, Chen, '20, OR'24



Find policy based on empirical MDP w/ slightly perturbed rewards

Optimal ℓ_∞ -based sample complexity

Theorem 4 (Li, Wei, Chi, Chen '20, OR'24)

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

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

with high prob., with sample complexity at most

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

Optimal ℓ_∞ -based sample complexity

Theorem 4 (Li, Wei, Chi, Chen '20, OR'24)

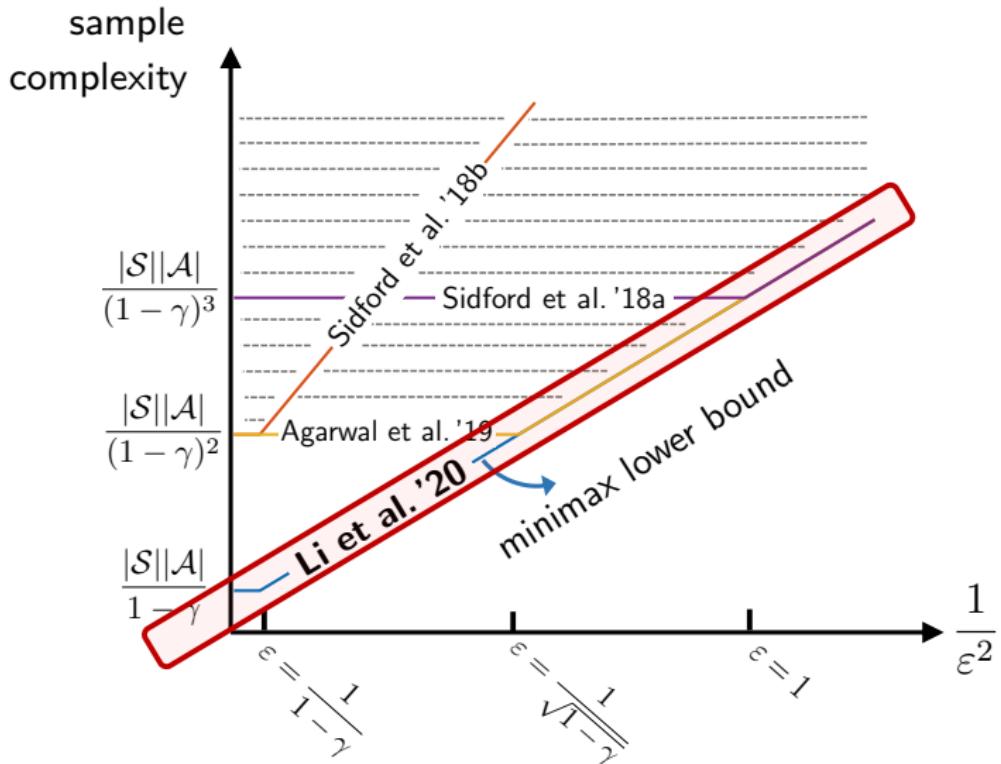
For any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the optimal policy $\widehat{\pi}_P^*$ of perturbed empirical MDP achieves

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

with high prob., with sample complexity at most

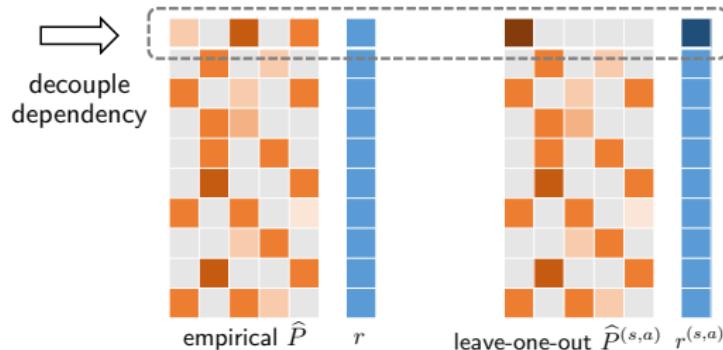
$$\tilde{O}\left(\frac{SA}{(1-\gamma)^3\varepsilon^2}\right)$$

- matches minimax lower bound: $\widetilde{\Omega}\left(\frac{SA}{(1-\gamma)^3\varepsilon^2}\right)$ Azar et al., 2013
- full ε -range: $\varepsilon \in (0, \frac{1}{1-\gamma}] \rightarrow$ no burn-in cost



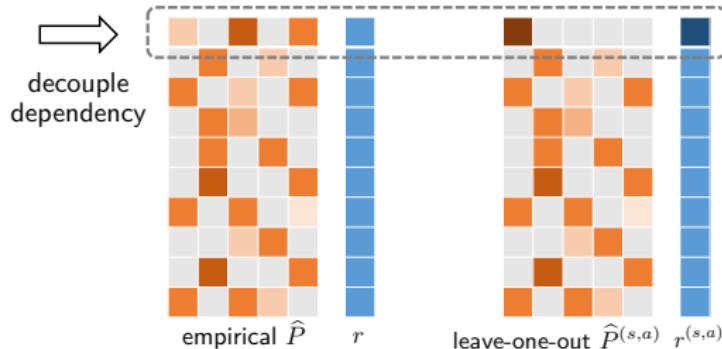
A glimpse of key analysis ideas

1. leave-one-out analysis: decouple statistical dependency by introducing auxiliary state-action absorbing MDPs by dropping randomness for each (s, a)



A glimpse of key analysis ideas

1. leave-one-out analysis: decouple statistical dependency by introducing auxiliary state-action absorbing MDPs by dropping randomness for each (s, a)

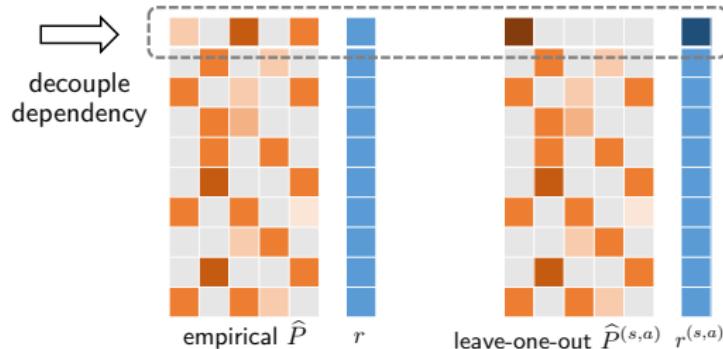


2. tie-breaking via random perturbation

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

A glimpse of key analysis ideas

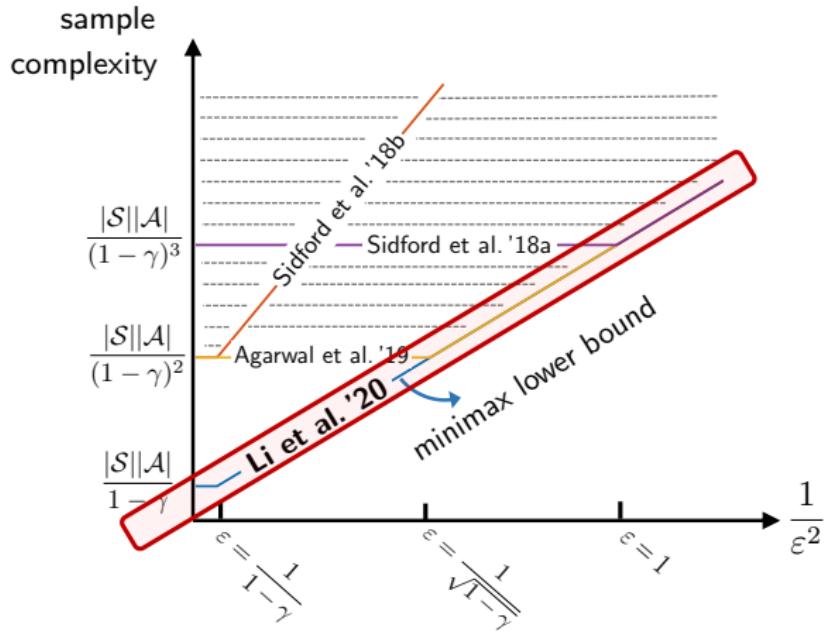
1. leave-one-out analysis: decouple statistical dependency by introducing auxiliary state-action absorbing MDPs by dropping randomness for each (s, a)



2. tie-breaking via random perturbation

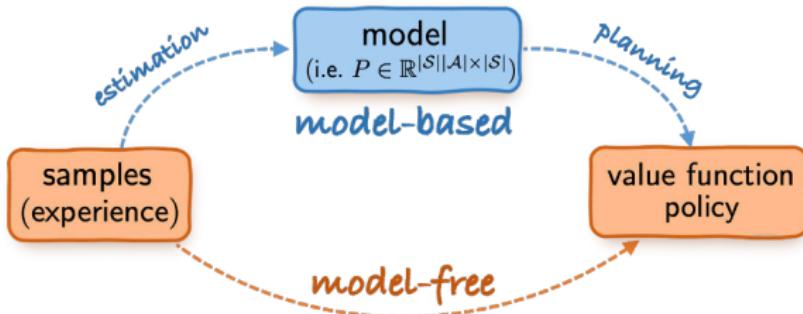
$$\forall s, \quad \hat{Q}^*(s, \hat{\pi}^*(s)) - \max_{a: a \neq \hat{\pi}^*(s)} \hat{Q}^*(s, a) > 0$$

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



Model based RL is minimax optimal under generative models
and does NOT suffer from a sample size barrier

Model-based vs. model-free RL



Model-based approach (“plug-in”)

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

Model-free / value-based approach

- learning w/o modeling & estimating environment explicitly
- memory-efficient, online, ...

Q-learning: a stochastic approximation algorithm



Chris Watkins



Peter Dayan

Stochastic approximation for solving the **Bellman equation**

Robbins & Monro, 1951

$$\mathcal{T}(Q) - Q = 0$$

where

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

Q-learning: a stochastic approximation algorithm



Chris Watkins



Peter Dayan

Stochastic approximation for solving Bellman equation $\mathcal{T}(Q) - Q = 0$

$$\underbrace{Q_{t+1}(s, a) = Q_t(s, a) + \eta_t(\mathcal{T}_t(Q_t)(s, a) - Q_t(s, a))}_{\text{sample transition } (s, a, s')}, \quad t \geq 0$$

Q-learning: a stochastic approximation algorithm



Chris Watkins



Peter Dayan

Stochastic approximation for solving Bellman equation $\mathcal{T}(Q) - Q = 0$

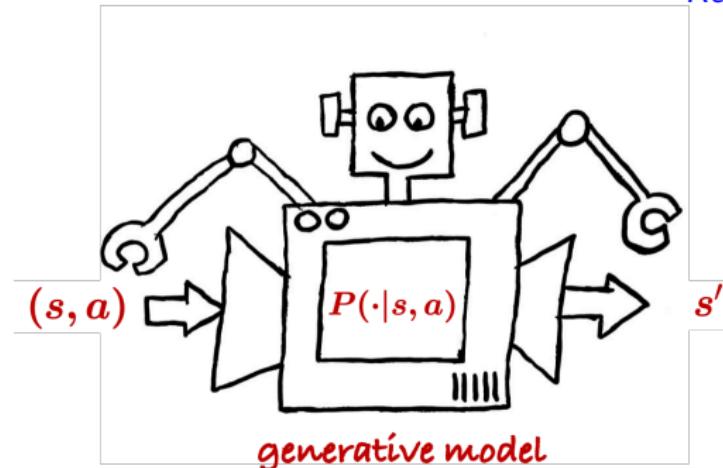
$$\underbrace{Q_{t+1}(s, a) = Q_t(s, a) + \eta_t (\mathcal{T}_t(Q_t)(s, a) - Q_t(s, a))}_{\text{sample transition } (s, a, s')} , \quad t \geq 0$$

$$\mathcal{T}_t(Q)(s, a) = r(s, a) + \gamma \max_{a'} Q(s', a')$$

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

A generative model / simulator

— Kearns, Singh, 1999



Each iteration, draw an independent sample (s, a, s') for given (s, a)

Synchronous Q-learning



Chris Watkins



Peter Dayan

for $t = 0, 1, \dots, T$

for each $(s, a) \in \mathcal{S} \times \mathcal{A}$

draw a sample (s, a, s') , run

$$Q_{t+1}(s, a) = (1 - \eta_t)Q_t(s, a) + \eta_t \left\{ r(s, a) + \gamma \max_{a'} Q_t(s', a') \right\}$$

synchronous: all state-action pairs are updated simultaneously

- total sample size: TSA

Sample complexity of synchronous Q-learning

Theorem 5 (Li, Cai, Chen, Wei, Chi '21, OR'24)

For any $0 < \varepsilon \leq 1$, synchronous Q-learning yields $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$ with high prob. and $\mathbb{E}[\|\hat{Q} - Q^*\|_\infty] \leq \varepsilon$, with sample size **at most**

$$\begin{cases} \tilde{O}\left(\frac{SA}{(1-\gamma)^4\varepsilon^2}\right) & \text{if } A \geq 2 \\ \tilde{O}\left(\frac{S}{(1-\gamma)^3\varepsilon^2}\right) & \text{if } A = 1 \end{cases} \quad (\text{TD learning})$$

Sample complexity of synchronous Q-learning

Theorem 5 (Li, Cai, Chen, Wei, Chi '21, OR'24)

For any $0 < \varepsilon \leq 1$, synchronous Q-learning yields $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$ with high prob. and $\mathbb{E}[\|\hat{Q} - Q^*\|_\infty] \leq \varepsilon$, with sample size **at most**

$$\begin{cases} \tilde{O}\left(\frac{SA}{(1-\gamma)^4\varepsilon^2}\right) & \text{if } A \geq 2 \\ \tilde{O}\left(\frac{S}{(1-\gamma)^3\varepsilon^2}\right) & \text{if } A = 1 \end{cases} \quad (\text{TD learning})$$

- Covers both *constant* and *rescaled linear* learning rates:

$$\eta_t \equiv \frac{1}{1 + \frac{c_1(1-\gamma)T}{\log^2 T}} \quad \text{or} \quad \eta_t = \frac{1}{1 + \frac{c_2(1-\gamma)t}{\log^2 T}}$$

Sample complexity of synchronous Q-learning

Theorem 5 (Li, Cai, Chen, Wei, Chi '21, OR'24)

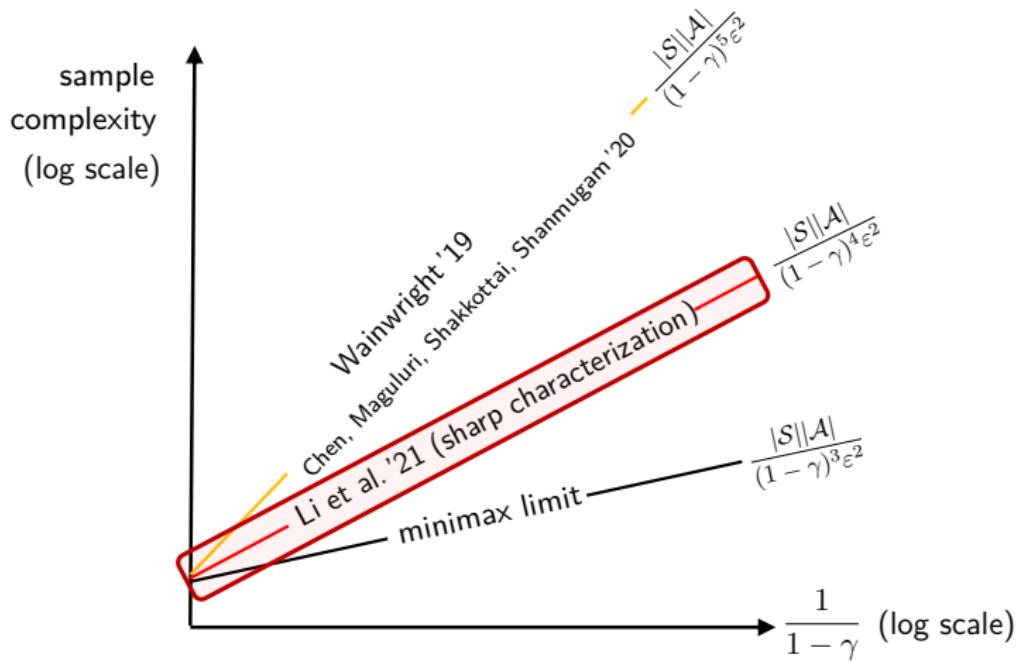
For any $0 < \varepsilon \leq 1$, synchronous Q-learning yields $\|\widehat{Q} - Q^*\|_\infty \leq \varepsilon$ with high prob. and $\mathbb{E}[\|\widehat{Q} - Q^*\|_\infty] \leq \varepsilon$, with sample size **at most**

$$\begin{cases} \tilde{O}\left(\frac{SA}{(1-\gamma)^4\varepsilon^2}\right) & \text{if } A \geq 2 \\ \tilde{O}\left(\frac{S}{(1-\gamma)^3\varepsilon^2}\right) & \text{if } A = 1 \end{cases} \quad (?)$$

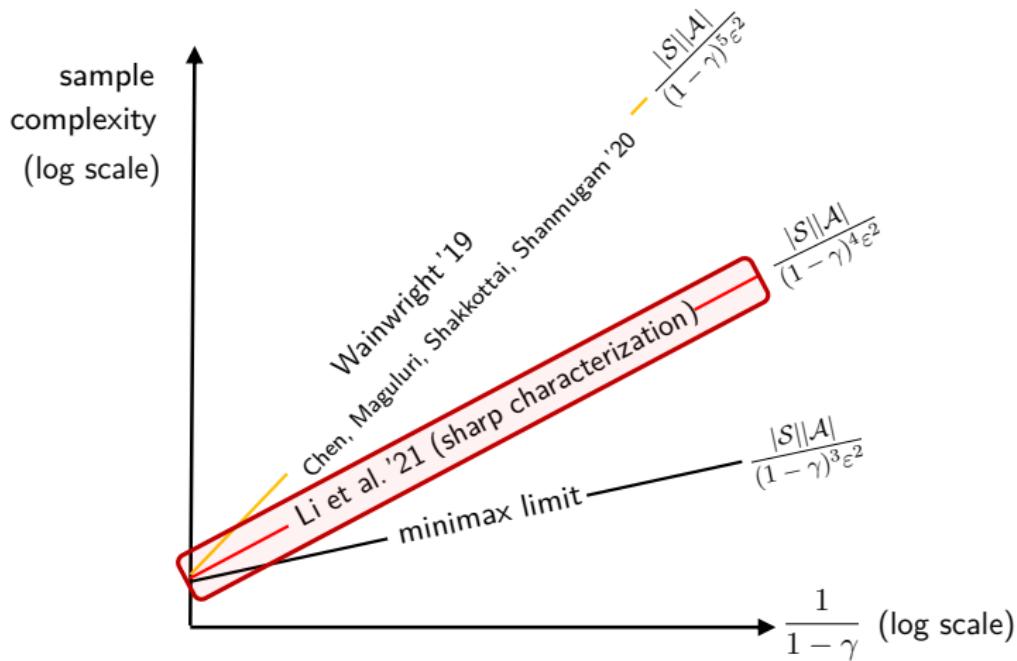
(minimax optimal)

other papers	sample complexity
Even-Dar & Mansour, 2003	$2^{\frac{1}{1-\gamma}} \frac{SA}{(1-\gamma)^4\varepsilon^2}$
Beck, Srikant, 2012	$\frac{S^2A^2}{(1-\gamma)^5\varepsilon^2}$
Wainwright, 2019	$\frac{SA}{(1-\gamma)^5\varepsilon^2}$
Chen, Maguluri, Shakkottai, Shanmugam, 2020	$\frac{SA}{(1-\gamma)^5\varepsilon^2}$

All this requires sample size at least $\frac{|S||A|}{(1-\gamma)^4 \varepsilon^2} (A \geq 2) \dots$



All this requires sample size at least $\frac{|S||A|}{(1-\gamma)^4 \varepsilon^2} (A \geq 2) \dots$



Question: Is Q-learning sub-optimal, or is it an analysis artifact?

Q-learning is NOT minimax optimal

Theorem 6 (Li, Cai, Chen, Wei, Chi '21, OR'24)

For any $0 < \varepsilon \leq 1$, there exists an MDP with $A \geq 2$ such that to achieve $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$, synchronous Q-learning needs *at least*

$$\tilde{\Omega}\left(\frac{SA}{(1-\gamma)^4\varepsilon^2}\right) \text{ samples}$$

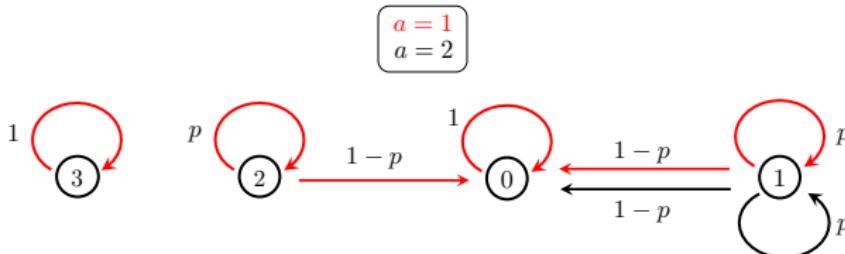
Q-learning is NOT minimax optimal

Theorem 6 (Li, Cai, Chen, Wei, Chi '21, OR'24)

For any $0 < \varepsilon \leq 1$, there exists an MDP with $A \geq 2$ such that to achieve $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$, synchronous Q-learning needs **at least**

$$\tilde{\Omega}\left(\frac{SA}{(1-\gamma)^4\varepsilon^2}\right) \text{ samples}$$

- Tight **algorithm-dependent** lower bound
- Holds for both constant and rescaled linear learning rates

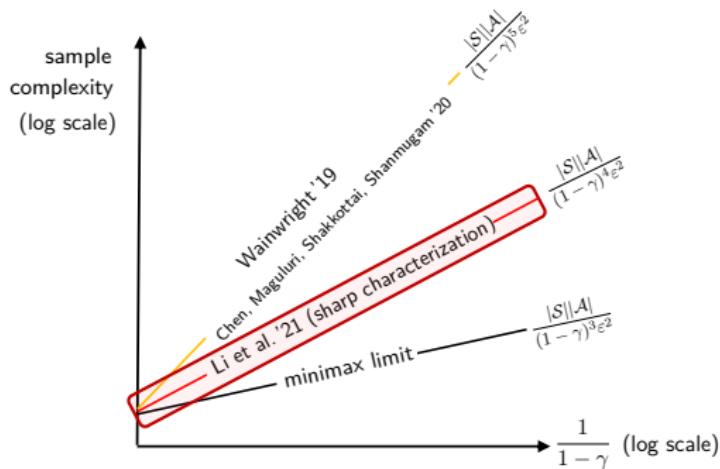


Q-learning is NOT minimax optimal

Theorem 6 (Li, Cai, Chen, Wei, Chi '21, OR'24)

For any $0 < \varepsilon \leq 1$, there exists an MDP with $A \geq 2$ such that to achieve $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$, synchronous Q-learning needs **at least**

$$\tilde{\Omega}\left(\frac{SA}{(1-\gamma)^4\varepsilon^2}\right) \text{ samples}$$



Why is Q-learning sub-optimal?

Over-estimation of Q-functions (Thrun & Schwartz '93; Hasselt '10)

- $\max_{a \in \mathcal{A}} \mathbb{E}[X(a)]$ tends to be over-estimated (high positive bias) when $\mathbb{E}[X(a)]$ is replaced by its empirical estimates using a small sample size
- often gets worse with a large number of actions (Hasselt, Guez, Silver '15)

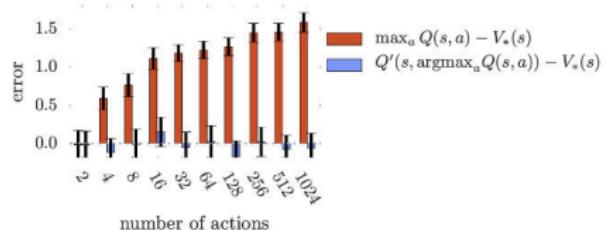


Figure 1: The orange bars show the bias in a single Q-learning update when the action values are $Q(s, a) = V_*(s) + \epsilon_a$ and the errors $\{\epsilon_a\}_{a=1}^m$ are independent standard normal random variables. The second set of action values Q' , used for the blue bars, was generated identically and independently. All bars are the average of 100 repetitions.

Why is Q-learning sub-optimal?

Over-estimation of Q-functions (Thrun & Schwartz '93; Hasselt '10)

- $\max_{a \in \mathcal{A}} \mathbb{E}[X(a)]$ tends to be over-estimated (high positive bias) when $\mathbb{E}[X(a)]$ is replaced by its empirical estimates using a small sample size
- often gets worse with a large number of actions (Hasselt, Guez, Silver '15)

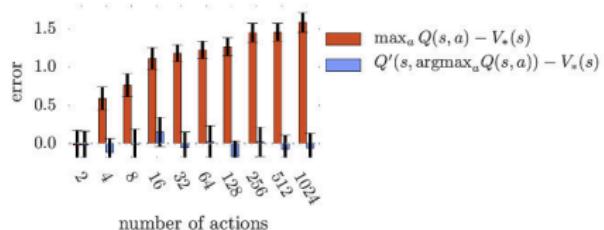


Figure 1: The orange bars show the bias in a single Q-learning update when the action values are $Q(s, a) = V_*(s) + \epsilon_a$ and the errors $\{\epsilon_a\}_{a=1}^m$ are independent standard normal random variables. The second set of action values Q' , used for the blue bars, was generated identically and independently. All bars are the average of 100 repetitions.

A provable improvement: Q-learning with variance reduction

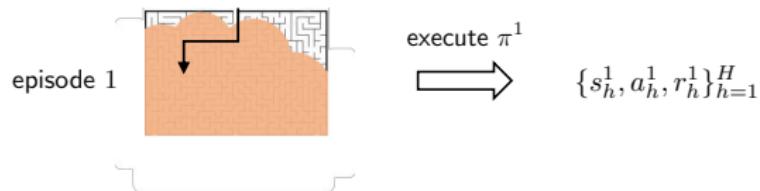
(Wainwright 2019)

Part 1

1. Basics: Markov decision processes
2. RL w/ a generative model (simulator)
3. Online RL
4. Offline RL

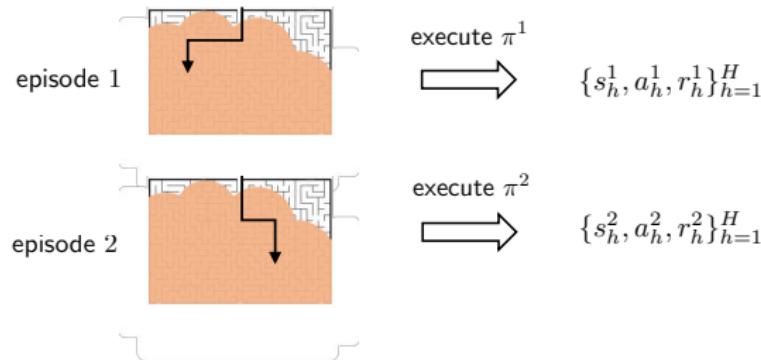
Online episodic RL

Sequentially execute MDP for K episodes, each consisting of H steps



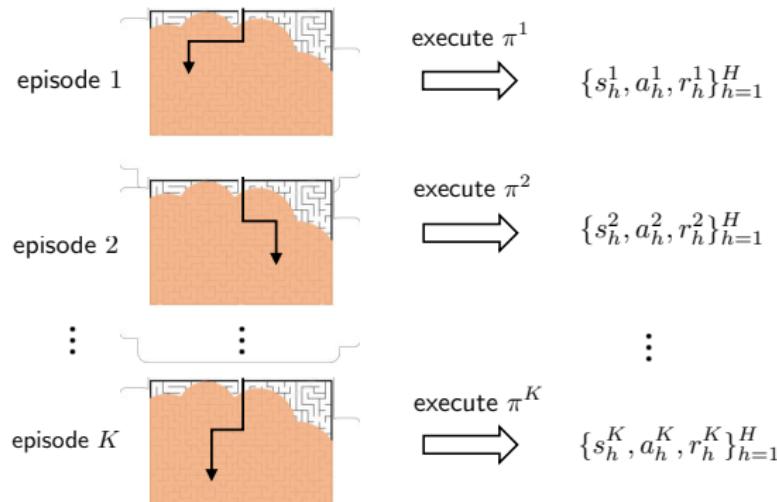
Online episodic RL

Sequentially execute MDP for K episodes, each consisting of H steps



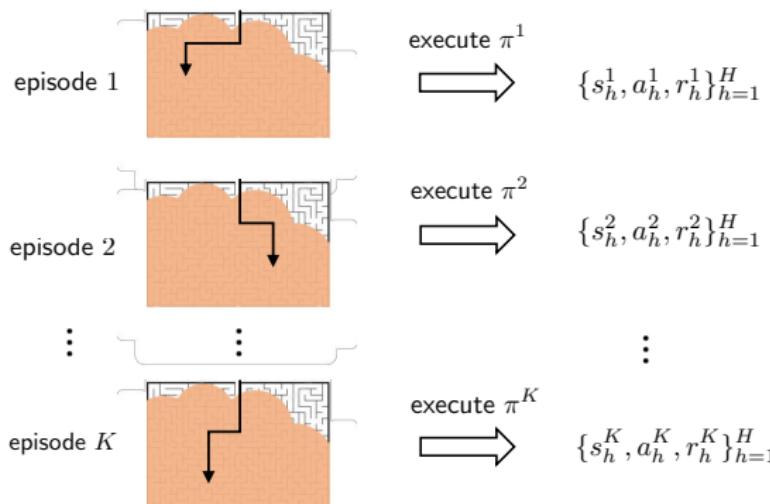
Online episodic RL

Sequentially execute MDP for K episodes, each consisting of H steps



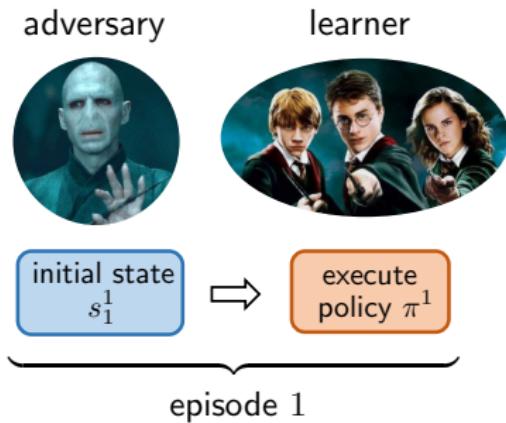
Online episodic RL

Sequentially execute MDP for K episodes, each consisting of H steps
— sample size: $T = KH$

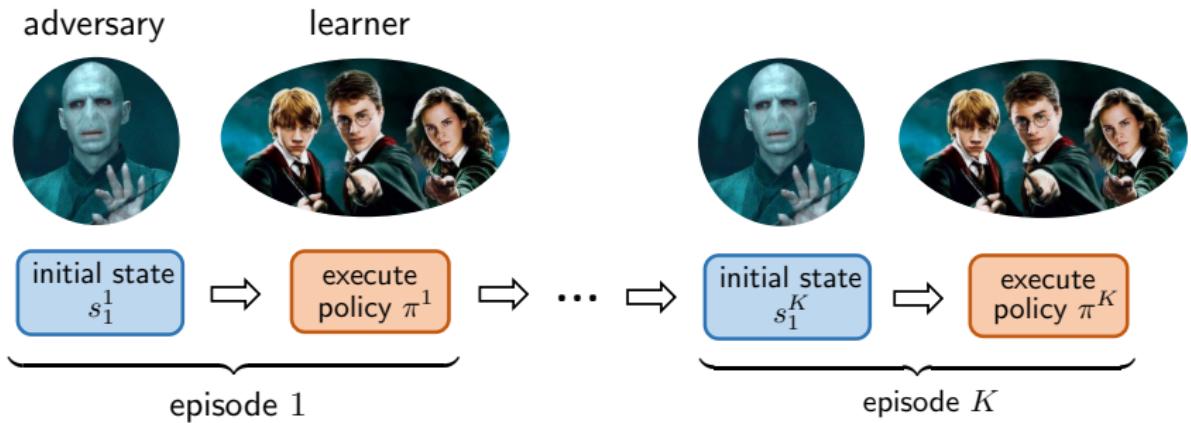


exploration (exploring unknowns) vs. **exploitation** (exploiting learned info)

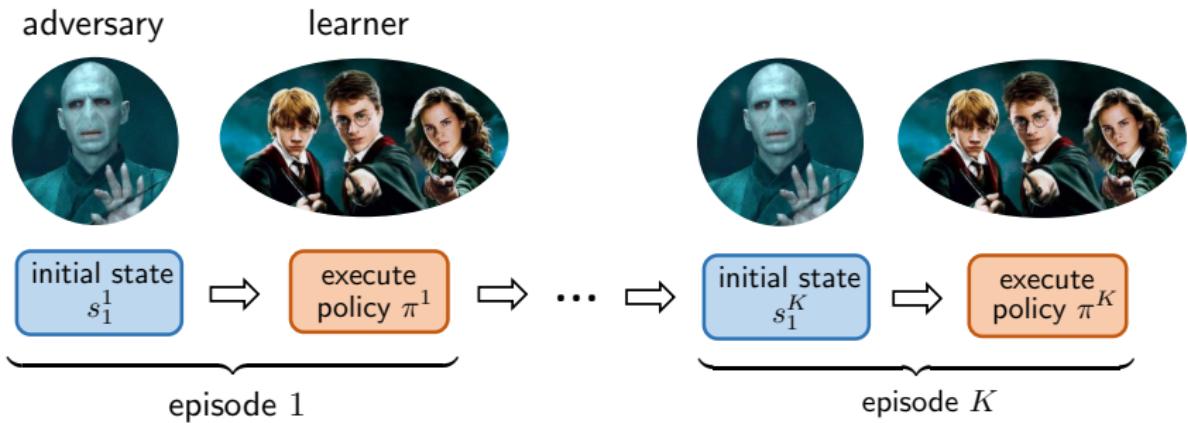
Regret: gap between learned policy & optimal policy



Regret: gap between learned policy & optimal policy



Regret: gap between learned policy & optimal policy



Performance metric: given initial states $\{s_1^k\}_{k=1}^K$, define

$$\text{Regret}(T) := \sum_{k=1}^K \left(V_1^\star(s_1^k) - V_1^{\pi^k}(s_1^k) \right)$$

Existing algorithms

- UCB-VI: [Azar et al, 2017](#)
- UBEV: [Dann et al, 2017](#)
- UCB-Q-Hoeffding: [Jin et al, 2018](#)
- UCB-Q-Bernstein: [Jin et al, 2018](#)
- UCB2-Q-Bernstein: [Bai et al, 2019](#)
- EULER: [Zanette et al, 2019](#)
- UCB-Q-Advantage: [Zhang et al, 2020](#)
- MVP: [Zhang et al, 2020](#)
- UCB-M-Q: [Menard et al, 2021](#)
- Q-EarlySettled-Advantage: [Li et al, 2021](#)
- (modified) MVP: [Zhang et al, 2024](#)

Lower bound

([Domingues et al, 2021](#))

$$\text{Regret}(T) \gtrsim \sqrt{H^2 SAT}$$

Existing algorithms

- UCB-VI: [Azar et al, 2017](#)
- UBEV: [Dann et al, 2017](#)
- UCB-Q-Hoeffding: [Jin et al, 2018](#)
- UCB-Q-Bernstein: [Jin et al, 2018](#)
- UCB2-Q-Bernstein: [Bai et al, 2019](#)
- EULER: [Zanette et al, 2019](#)
- UCB-Q-Advantage: [Zhang et al, 2020](#)
- MVP: [Zhang et al, 2020](#)
- UCB-M-Q: [Menard et al, 2021](#)
- Q-EarlySettled-Advantage: [Li et al, 2021](#)
- (modified) MVP: [Zhang et al, 2024](#)

Lower bound

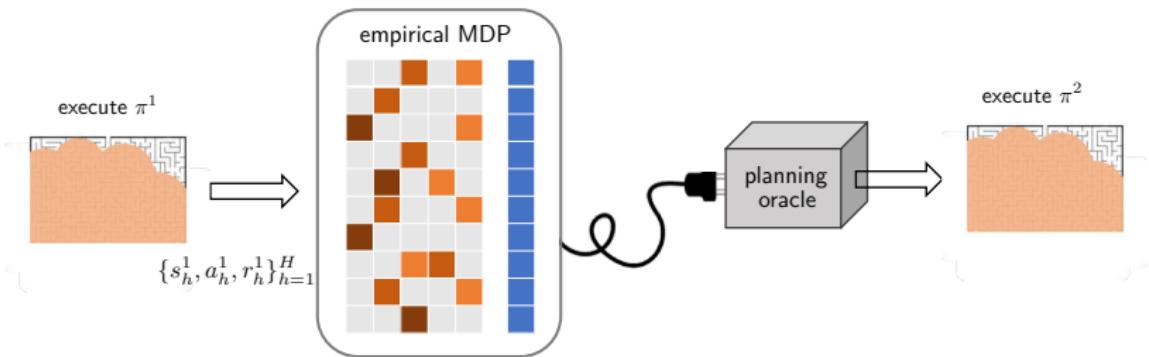
([Domingues et al, 2021](#))

$$\text{Regret}(T) \gtrsim \sqrt{H^2 SAT}$$

Which online RL algorithms achieve near-minimal regret?

Model-based online RL with UCB exploration

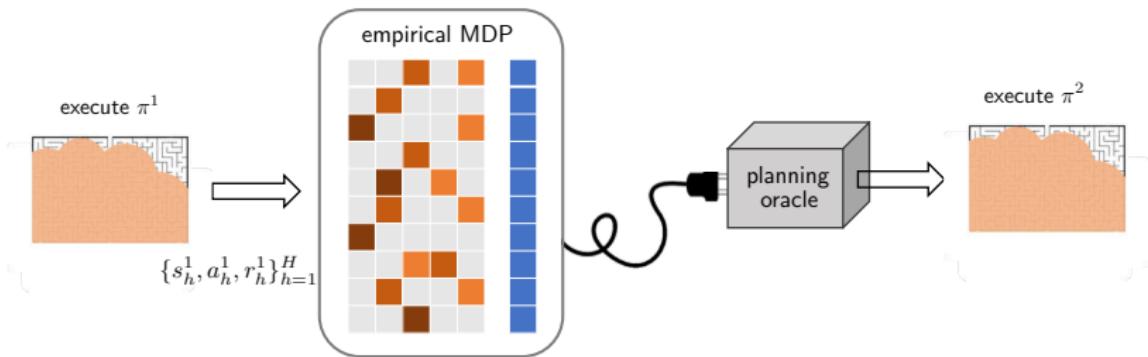
Model-based approach for online RL



repeat:

- use collected data to estimate transition probabilities
- apply planning to the estimated model to derive a new policy for sampling in the next episode

Model-based approach for online RL



repeat:

- use collected data to estimate transition probabilities
- apply planning to the estimated model to derive a new policy for sampling in the next episode

How to balance exploration and exploitation in this framework?



T. L. Lai

H. Robbins

Optimism in the face of uncertainty:

- explores based on the best optimistic estimates associated with the actions!
- a common framework: utilize upper confidence bounds (UCB)
accounts for estimates + uncertainty level



T. L. Lai



H. Robbins

Optimism in the face of uncertainty:

- explores based on the best optimistic estimates associated with the actions!
- a common framework: utilize upper confidence bounds (UCB)
accounts for estimates + uncertainty level

Optimistic model-based approach: incorporates **UCB** framework into model-based approach

UCB-VI (Azar et al. '17)

For each episode:

1. Backtrack $h = H, H - 1, \dots, 1$: run **value iteration**

$$Q_h(s_h, a_h) \leftarrow r_h(s_h, a_h) + \underbrace{\hat{P}_{h, s_h, a_h}}_{\text{model estimate}} V_{h+1}$$

$$V_h(s_h) \leftarrow \max_{a \in \mathcal{A}} Q_h(s_h, a)$$

UCB-VI (Azar et al. '17)

For each episode:

1. Backtrack $h = H, H - 1, \dots, 1$: run **optimistic value iteration**

$$Q_h(s_h, a_h) \leftarrow r_h(s_h, a_h) + \underbrace{\hat{P}_{h, s_h, a_h}}_{\text{model estimate}} V_{h+1} + \underbrace{b_h(s_h, a_h)}_{\text{bonus (upper confidence width)}}$$
$$V_h(s_h) \leftarrow \max_{a \in \mathcal{A}} Q_h(s_h, a)$$

UCB-VI (Azar et al. '17)

For each episode:

1. Backtrack $h = H, H - 1, \dots, 1$: run **optimistic value iteration**

$$Q_h(s_h, a_h) \leftarrow r_h(s_h, a_h) + \underbrace{\hat{P}_{h, s_h, a_h}}_{\text{model estimate}} V_{h+1} + \underbrace{b_h(s_h, a_h)}_{\text{bonus (upper confidence width)}}$$

$$V_h(s_h) \leftarrow \max_{a \in \mathcal{A}} Q_h(s_h, a)$$

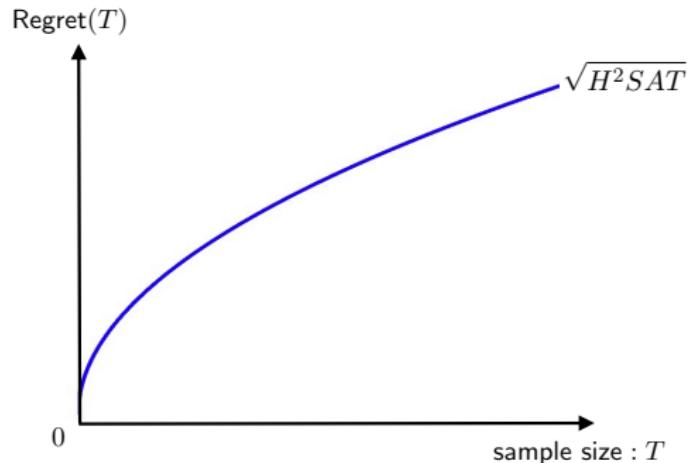
2. Forward $h = 1, \dots, H$: take actions according to **greedy policy**

$$\pi_h(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} Q_h(s, a)$$

to sample a new episode $\{s_h, a_h, r_h\}_{h=1}^H$

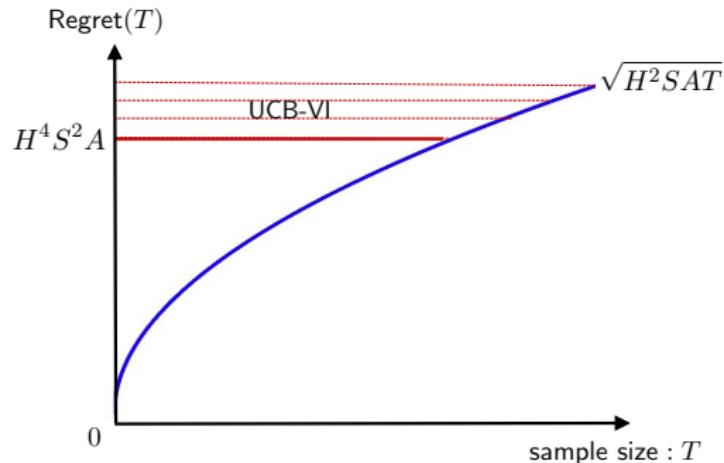
UCB-VI is asymptotically regret-optimal

— Azar, Osband, Munos, 2017



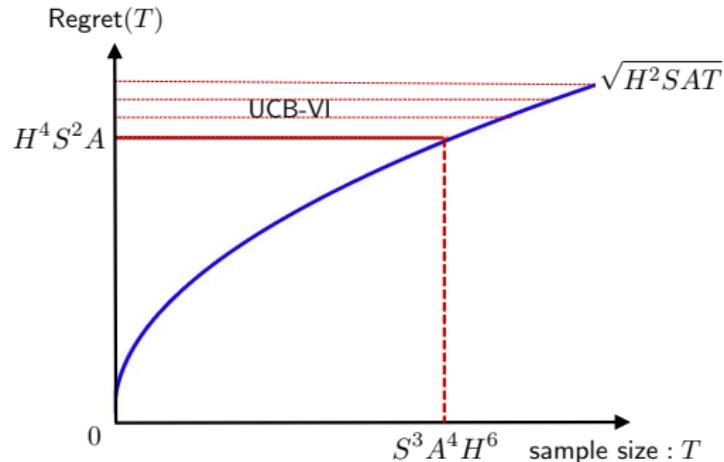
UCB-VI is asymptotically regret-optimal

— Azar, Osband, Munos, 2017



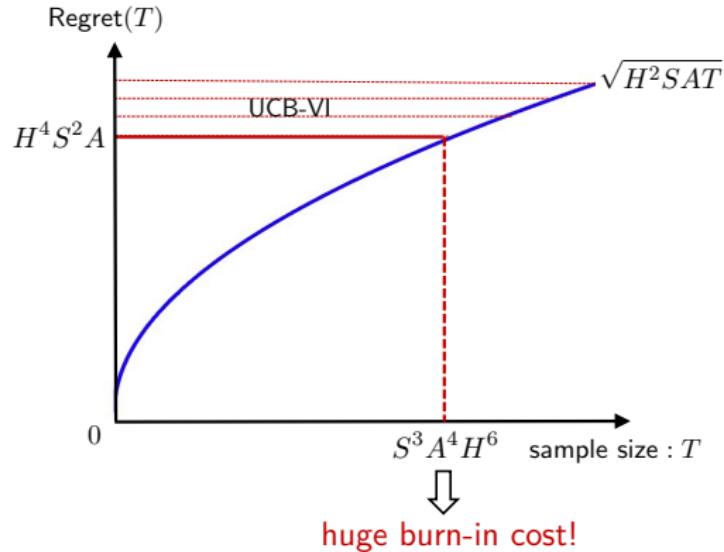
UCB-VI is asymptotically regret-optimal

— Azar, Osband, Munos, 2017



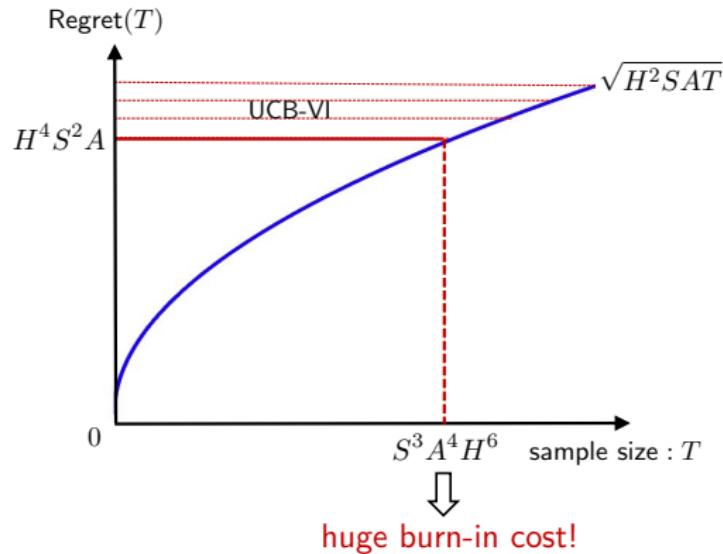
UCB-VI is asymptotically regret-optimal

— Azar, Osband, Munos, 2017



UCB-VI is asymptotically regret-optimal

— Azar, Osband, Munos, 2017



Issues: large burn-in cost

Other asymptotically regret-optimal algorithms

Algorithm	Regret upper bound	Range of K that attains optimal regret
UCBVI (Azar et al, 2017)	$\sqrt{SAH^2T} + S^2 AH^3$	$[S^3 AH^3, \infty)$
ORLC (Dann et al, 2019)	$\sqrt{SAH^2T} + S^2 AH^4$	$[S^3 AH^5, \infty)$
EULER (Zanette et al, 2019)	$\sqrt{SAH^2T} + S^{3/2} AH^3(\sqrt{S} + \sqrt{H})$	$[S^2 AH^3(\sqrt{S} + \sqrt{H}), \infty)$
UCB-Adv (Zhang et al, 2020)	$\sqrt{SAH^2T} + S^2 A^{3/2} H^{33/4} K^{1/4}$	$[S^6 A^4 H^{27}, \infty)$
MVP (Zhang et al, 2020)	$\sqrt{SAH^2T} + S^2 AH^2$	$[S^3 AH, \infty)$
UCB-M-Q (Menard et al, 2021)	$\sqrt{SAH^2T} + SAH^4$	$[SAH^5, \infty)$
Q-Earlysettled-Adv (Li et al, 2021)	$\sqrt{SAH^2T} + SAH^6$	$[SAH^9, \infty)$

Other asymptotically regret-optimal algorithms

Algorithm	Regret upper bound	Range of K that attains optimal regret
UCBVI (Azar et al, 2017)	$\sqrt{SAH^2T} + S^2 AH^3$	$[S^3 AH^3, \infty)$
ORLC (Dann et al, 2019)	$\sqrt{SAH^2T} + S^2 AH^4$	$[S^3 AH^5, \infty)$
EULER (Zanette et al, 2019)	$\sqrt{SAH^2T} + S^{3/2} AH^3(\sqrt{S} + \sqrt{H})$	$[S^2 AH^3(\sqrt{S} + \sqrt{H}), \infty)$
UCB-Adv (Zhang et al, 2020)	$\sqrt{SAH^2T} + S^2 A^{3/2} H^{33/4} K^{1/4}$	$[S^6 A^4 H^{27}, \infty)$
MVP (Zhang et al, 2020)	$\sqrt{SAH^2T} + S^2 AH^2$	$[S^3 AH, \infty)$
UCB-M-Q (Menard et al, 2021)	$\sqrt{SAH^2T} + SAH^4$	$[SAH^5, \infty)$
Q-Earlysettled-Adv (Li et al, 2021)	$\sqrt{SAH^2T} + SAH^6$	$[SAH^9, \infty)$

Can we find a regre-optimal algorithm with no burn-in cost?

Monotonic Value Propagation

UCB-VI with **doubling update rules** and **variance-aware bonus**

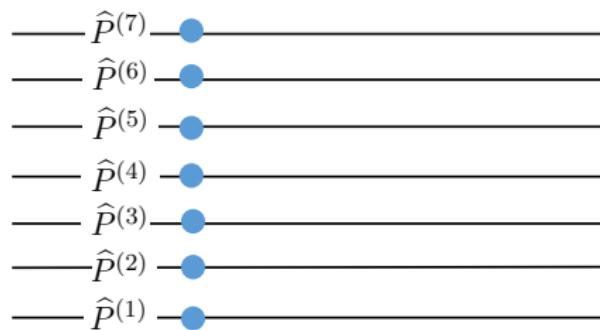
- (s, a, h) is updated only when visited the $\{1, 3, 7, 15, \dots\}$ -th time

Monotonic Value Propagation

UCB-VI with **doubling update rules** and **variance-aware bonus**

- (s, a, h) is updated only when visited the $\{1, 3, 7, 15, \dots\}$ -th time

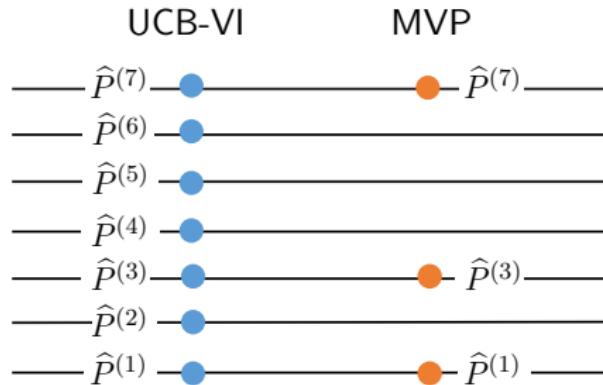
UCB-VI



Monotonic Value Propagation

UCB-VI with **doubling update rules** and **variance-aware bonus**

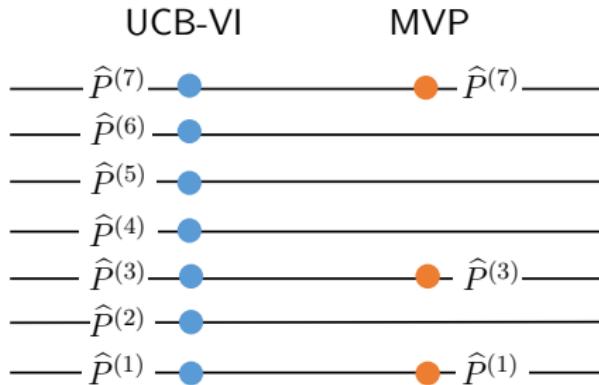
- (s, a, h) is updated only when visited the $\{1, 3, 7, 15, \dots\}$ -th time



Monotonic Value Propagation

UCB-VI with **doubling update rules** and **variance-aware bonus**

- (s, a, h) is updated only when visited the $\{1, 3, 7, 15, \dots\}$ -th time

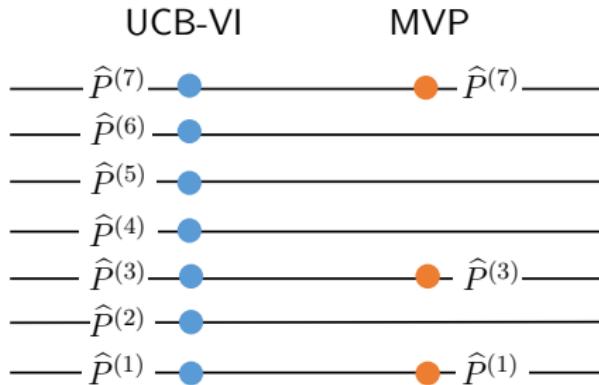


- visitation counts change much less frequently
→ reduces covering number dramatically

Monotonic Value Propagation

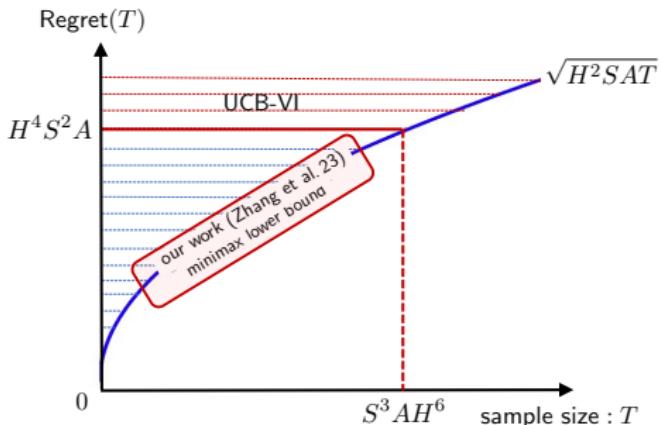
UCB-VI with **doubling update rules** and **variance-aware bonus**

- (s, a, h) is updated only when visited the $\{1, 3, 7, 15, \dots\}$ -th time



- visitation counts change much less frequently
→ reduces covering number dramatically
- data-driven bonus terms (chosen based on empirical variances)

Regret-optimal algorithm w/o burn-in cost

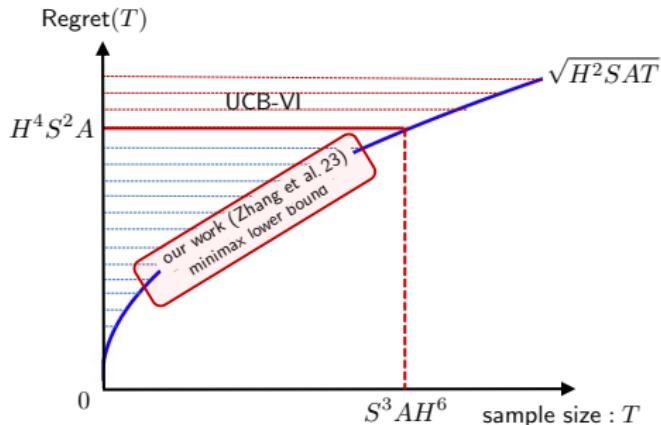


Theorem 7 (Zhang, Chen, Lee, Du '24)

The model-based algorithm Monotonic Value Propagation achieves

$$\text{Regret}(T) \lesssim \tilde{O}(\sqrt{H^2 SAT})$$

Regret-optimal algorithm w/o burn-in cost



Theorem 7 (Zhang, Chen, Lee, Du '24)

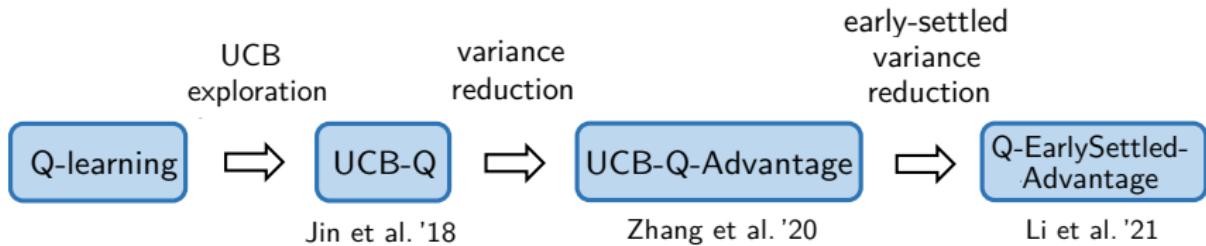
The model-based algorithm Monotonic Value Propagation achieves

$$\text{Regret}(T) \lesssim \tilde{O}(\sqrt{H^2 SAT})$$

- the only algorithm so far that is regret-optimal w/o burn-ins

Which model-free algorithms are sample-efficient for online RL?

Which model-free algorithms are sample-efficient for online RL?



Q-learning with UCB exploration (Jin et al., 2018)

$$Q_h(s_h, a_h) \leftarrow \underbrace{(1 - \eta_k)Q_h(s_h, a_h) + \eta_k \mathcal{T}_k(Q_{h+1})(s_h, a_h)}_{\text{classical Q-learning}} + \eta_k \underbrace{b_h(s_h, a_h)}_{\text{exploration bonus}}$$

Q-learning with UCB exploration (Jin et al., 2018)

$$Q_h(s_h, a_h) \leftarrow \underbrace{(1 - \eta_k)Q_h(s_h, a_h) + \eta_k \mathcal{T}_k(Q_{h+1})(s_h, a_h)}_{\text{classical Q-learning}} + \eta_k \underbrace{b_h(s_h, a_h)}_{\text{exploration bonus}}$$

- $b_h(s, a)$: upper confidence bound
 - *optimism in the face of uncertainty*
- inspired by UCB bandit algorithm (Lai, Robbins '85)

Q-learning with UCB exploration (Jin et al., 2018)

$$Q_h(s_h, a_h) \leftarrow \underbrace{(1 - \eta_k)Q_h(s_h, a_h) + \eta_k \mathcal{T}_k(Q_{h+1})(s_h, a_h)}_{\text{classical Q-learning}} + \eta_k \underbrace{b_h(s_h, a_h)}_{\text{exploration bonus}}$$

- $b_h(s, a)$: upper confidence bound
 - *optimism in the face of uncertainty*
- inspired by UCB bandit algorithm (Lai, Robbins '85)

Regret(T) $\lesssim \sqrt{H^3 SAT}$ \implies sub-optimal by a factor of \sqrt{H}

Q-learning with UCB exploration (Jin et al., 2018)

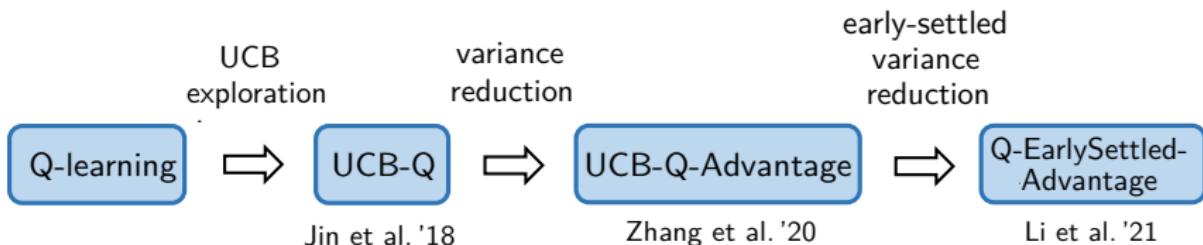
$$Q_h(s_h, a_h) \leftarrow \underbrace{(1 - \eta_k)Q_h(s_h, a_h) + \eta_k \mathcal{T}_k(Q_{h+1})(s_h, a_h)}_{\text{classical Q-learning}} + \eta_k \underbrace{b_h(s_h, a_h)}_{\text{exploration bonus}}$$

- $b_h(s, a)$: upper confidence bound
 - *optimism in the face of uncertainty*
- inspired by UCB bandit algorithm (Lai, Robbins '85)

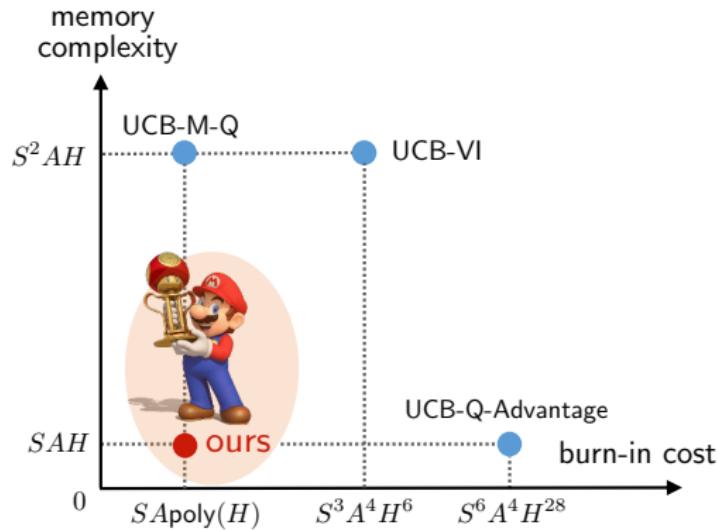
Regret(T) $\lesssim \sqrt{H^3 SAT}$ \implies sub-optimal by a factor of \sqrt{H}

Issue: large variability in stochastic update rules

Further improvement

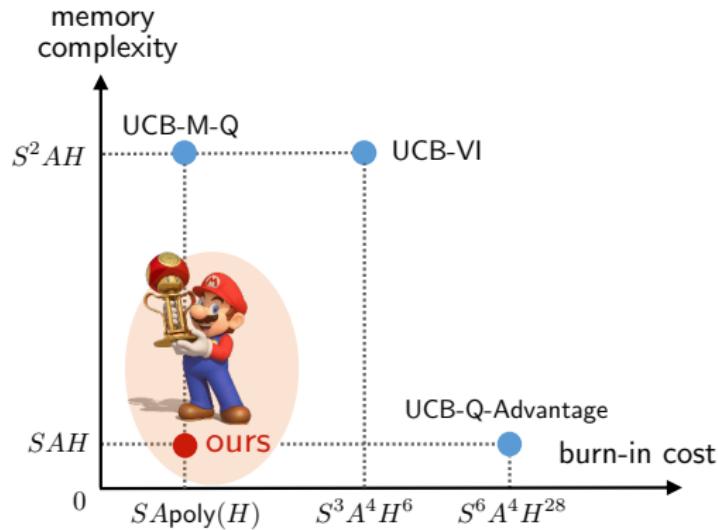


- UCB-Q-Advantage: use variance reduction to achieve near-optimal regret, but with large burn-in cost;
- Q-EarlySettled-Advantage: stop updating the reference as soon as possible to reduce burn-in cost.



Model-free algorithms can simultaneously achieve

- (1) regret optimality; (2) **low** burn-in cost; (3) memory efficiency



Model-free algorithms can simultaneously achieve

- (1) regret optimality; (2) **low** burn-in cost; (3) memory efficiency

Part 1

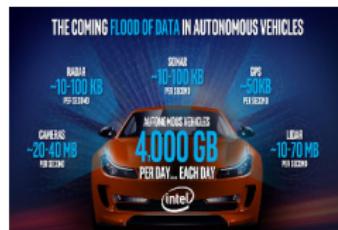
1. Basics: Markov decision processes
2. RL w/ a generative model (simulator)
3. Online RL
4. Offline RL

Offline/batch RL

- Collecting new data might be costly, unsafe, unethical, or time-consuming



medical records



data of self-driving



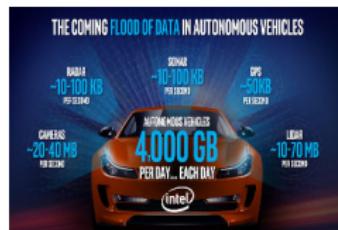
clicking times of ads

Offline/batch RL

- Collecting new data might be costly, unsafe, unethical, or time-consuming
- But we have already stored tons of historical data



medical records



data of self-driving



clicking times of ads

Offline/batch RL

- Collecting new data might be costly, unsafe, unethical, or time-consuming
- But we have already stored tons of historical data



medical records



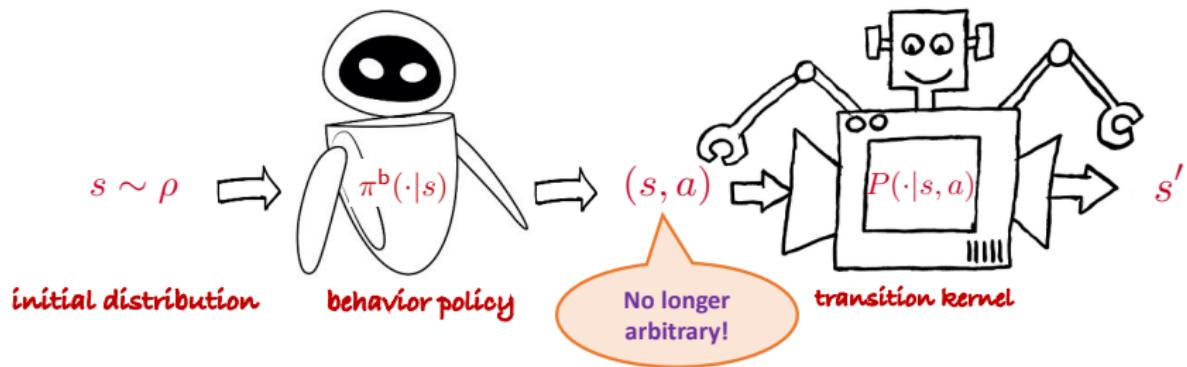
data of self-driving



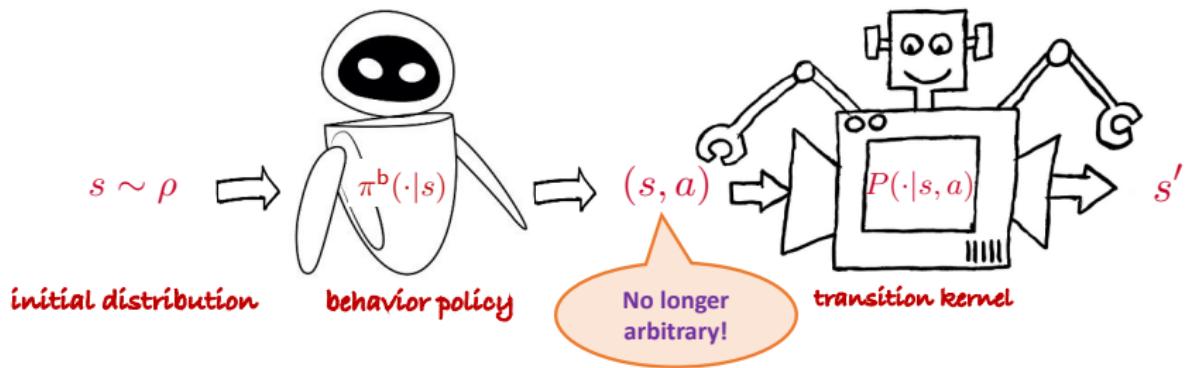
clicking times of ads

Question: can we learn based solely on historical data w/o active exploration?

A mathematical model of offline data



A mathematical model of offline data

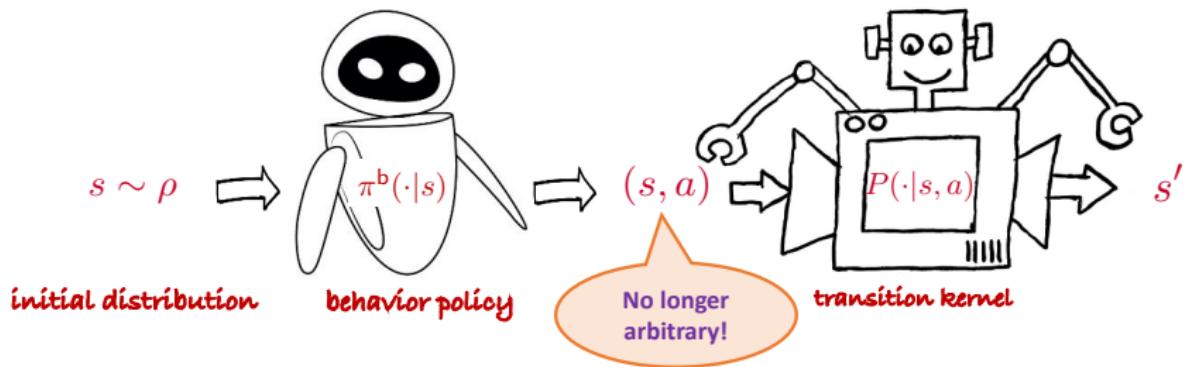


historical dataset $\mathcal{D} = \{(s^{(i)}, a^{(i)}, s'^{(i)})\}$: N independent copies of

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

- ρ : initial state distribution; π^b : behavior policy

A mathematical model of offline data



Goal: given a target accuracy level $\varepsilon \in (0, H]$, find $\hat{\pi}$ s.t.

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

— *in a sample-efficient manner*

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 (Rashidinejad et al. '21)

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

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

Single-policy concentrability coefficient (Rashidinejad et al. '21)

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

- captures distributional shift

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

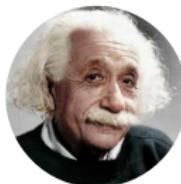
Single-policy concentrability coefficient (Rashidinejad et al. '21)

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

- captures distributional shift

$$C^* = O(1)$$

large C^*



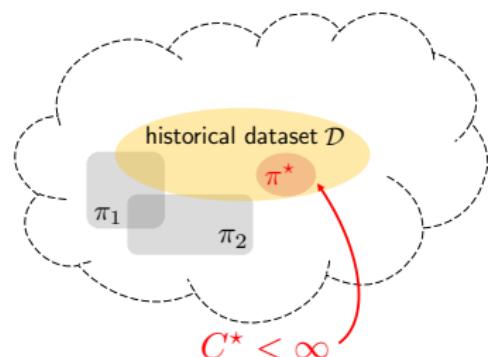
expert data

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

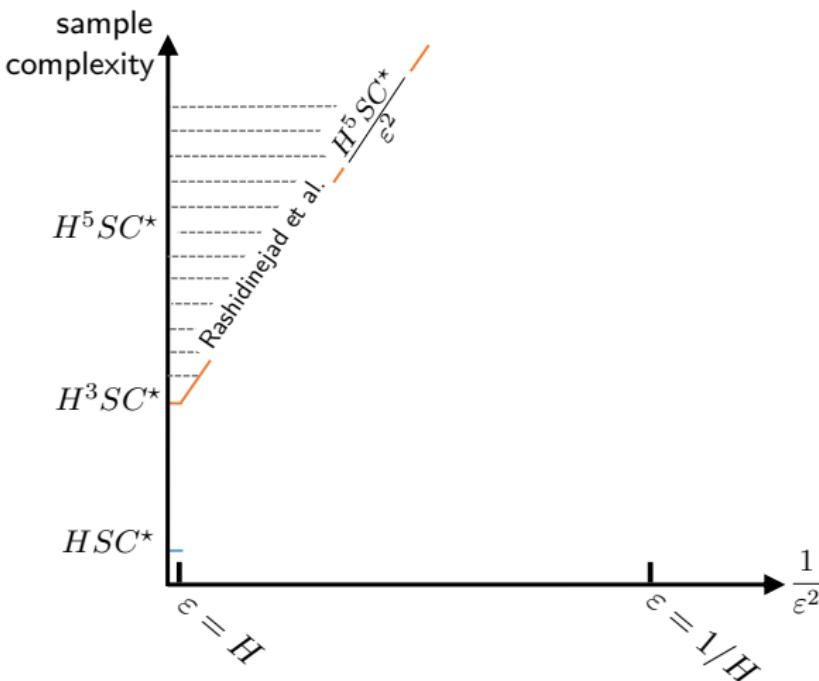
Single-policy concentrability coefficient (Rashidinejad et al. '21)

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

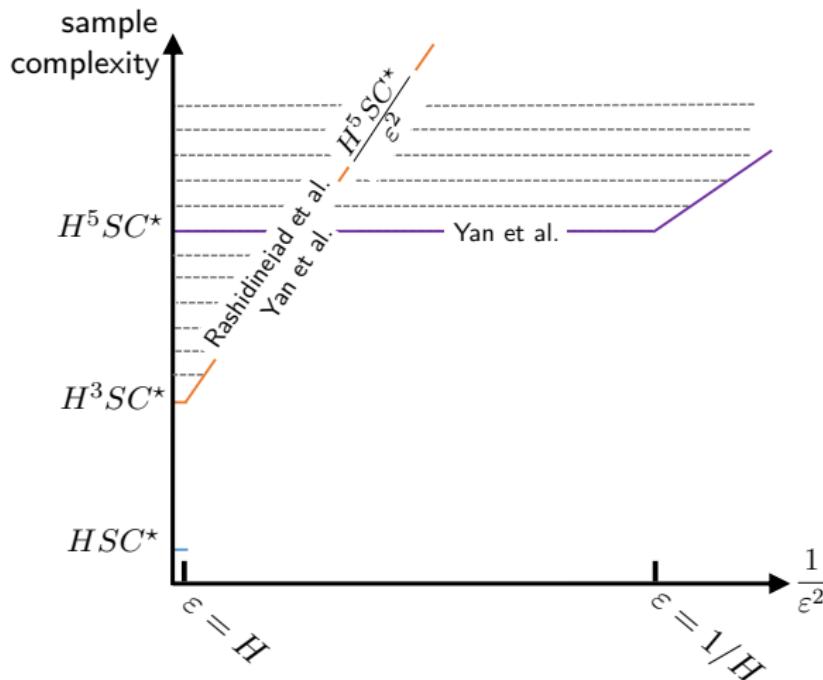
- captures distributional shift
- allows for partial coverage
 - as long as it covers the part reachable by π^*



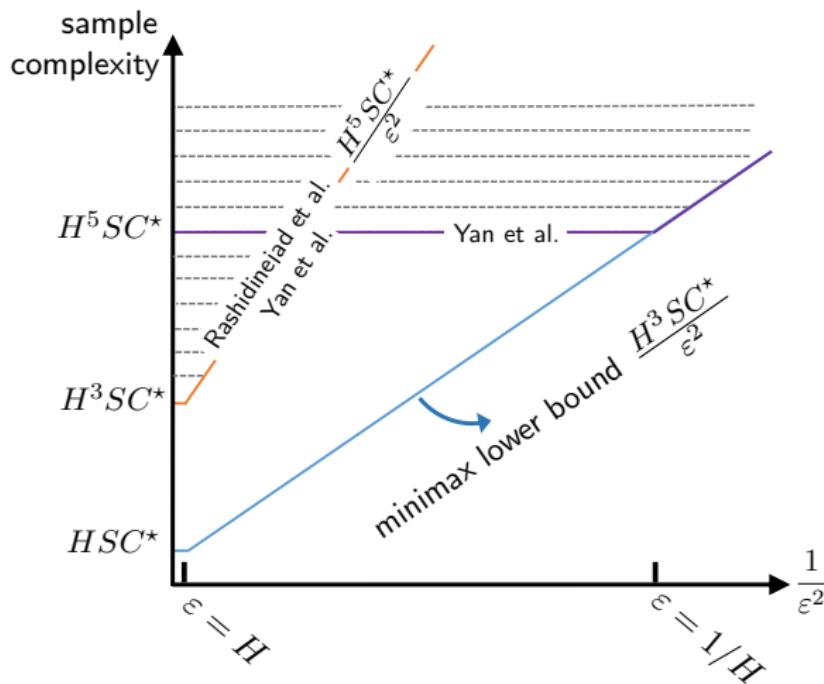
Prior art: sample complexity bounds



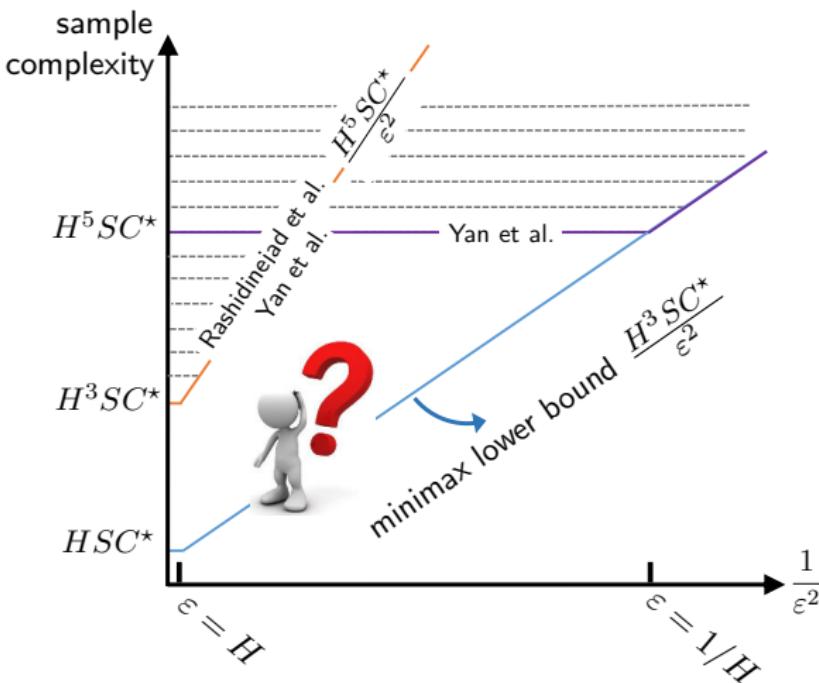
Prior art: sample complexity bounds



Prior art: sample complexity bounds

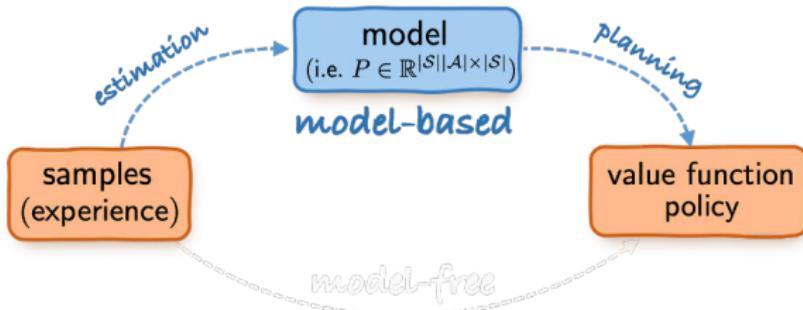


Prior art: sample complexity bounds

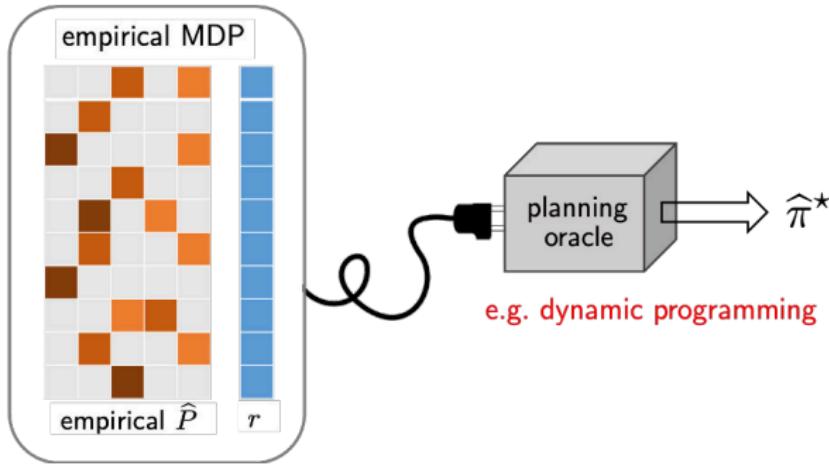


Can we close the gap between upper & lower bounds?

Model-based (“plug-in”) approach?



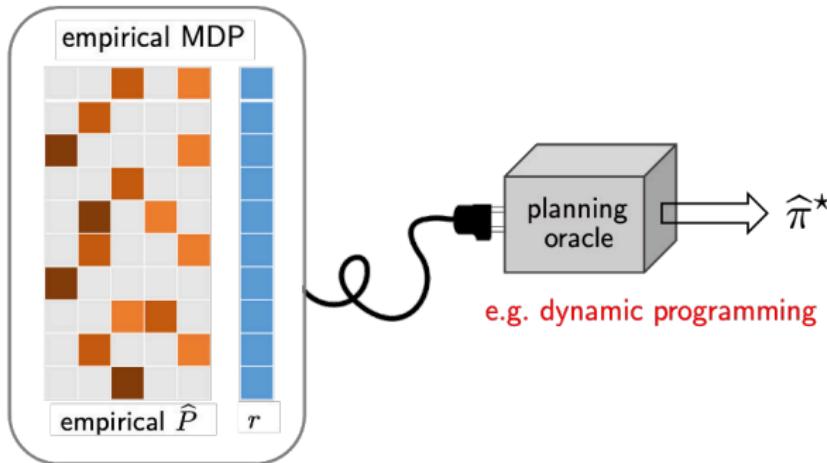
Model-based (“plug-in”) approach?



1. construct empirical model \hat{P} :

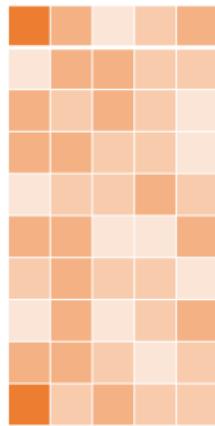
$$\hat{P}(s' | s, a) = \underbrace{\frac{1}{N} \sum_{i=1}^N \mathbf{1}\{s'^{(i)} = s'\}}_{\text{empirical frequency}}$$

Model-based (“plug-in”) approach?

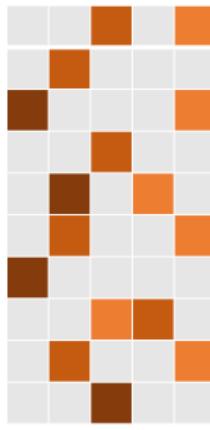


1. construct empirical model \hat{P}
2. planning (e.g. value iteration) based on empirical MDP

Issues & challenges in the sample-starved regime



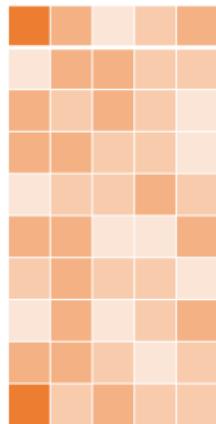
truth: $P \in \mathbb{R}^{SA \times S}$



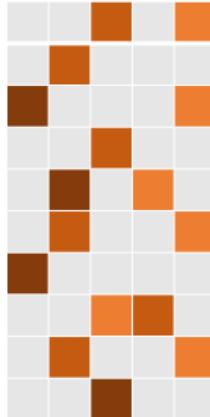
empirical \hat{P} (simulator)

- can't recover P faithfully if sample size $\ll S^2 A$

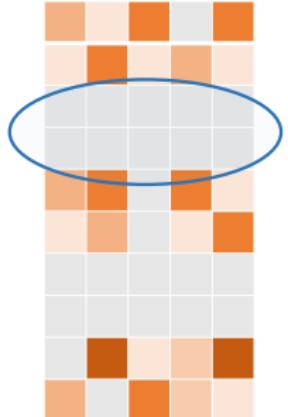
Issues & challenges in the sample-starved regime



truth: $P \in \mathbb{R}^{SA \times S}$



empirical \hat{P} (simulator)



empirical \hat{P} (offline)

- can't recover P faithfully if sample size $\ll S^2 A$
- (possibly) insufficient coverage under offline data

Key idea: pessimism in the face of uncertainty

— *Jin et al, 2020, Rashidinejad et al, 2021, Xie et al, 2021*



online

upper confidence bounds

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

Key idea: pessimism in the face of uncertainty

— Jin et al, 2020, Rashidinejad et al, 2021, Xie et al, 2021

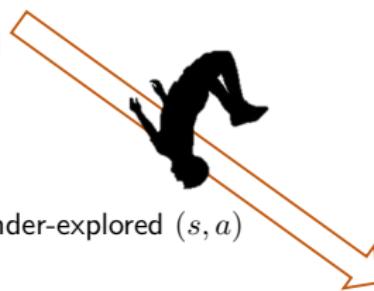


upper confidence bounds

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



online



offline



lower confidence bounds

— stay cautious about under-explored (s, a)

Key idea: pessimism in the face of uncertainty

— Jin et al, 2020, Rashidinejad et al, 2021, Xie et al, 2021

1. build empirical model \hat{P}
2. (**value iteration**) repeat: for all (s, a)

$$\hat{Q}(s, a) \leftarrow \max \left\{ r(s, a) + \gamma \langle \hat{P}(\cdot | s, a), \hat{V} \rangle, 0 \right\}$$

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

Key idea: pessimism in the face of uncertainty

— Jin et al, 2020, Rashidinejad et al, 2021, Xie et al, 2021

Penalize those poorly visited $(s, a) \dots$

1. build empirical model \hat{P}
2. (**pessimistic value iteration**) repeat: for all (s, a)

$$\hat{Q}(s, a) \leftarrow \max \left\{ r(s, a) + \gamma \langle \hat{P}(\cdot | s, a), \hat{V} \rangle - \underbrace{b(s, a; \hat{V})}_{\text{uncertainty penalty}}, 0 \right\}$$

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

Key idea: pessimism in the face of uncertainty

— Jin et al, 2020, Rashidinejad et al, 2021, Xie et al, 2021

Penalize those poorly visited $(s, a) \dots$

1. build empirical model \hat{P}
2. **(pessimistic value iteration)** repeat: for all (s, a)

$$\hat{Q}(s, a) \leftarrow \max \left\{ r(s, a) + \gamma \langle \hat{P}(\cdot | s, a), \hat{V} \rangle - \underbrace{b(s, a; \hat{V})}_{\text{uncertainty penalty}}, 0 \right\}$$

compared w/ Rashidinejad et al, 2021

- sample-reuse across iterations
- Bernstein-style penalty

Sample complexity of model-based offline RL

Theorem 8 (Li, Shi, Chen, Chi, Wei '24)

For any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the policy $\widehat{\pi}$ returned by VI-LCB using a Bernstein-style penalty term achieves

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

with high prob., with sample complexity at most

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

Sample complexity of model-based offline RL

Theorem 8 (Li, Shi, Chen, Chi, Wei '24)

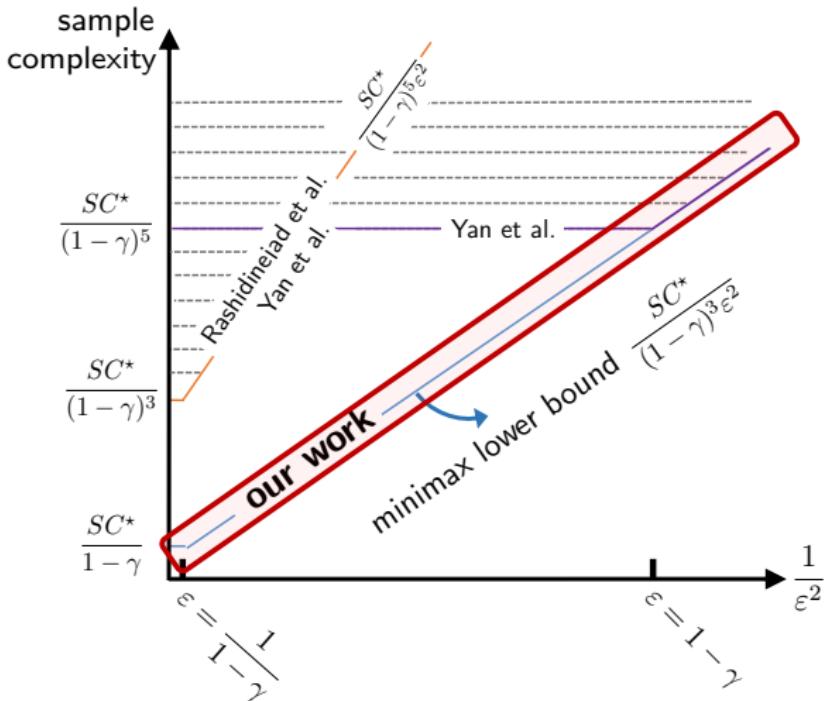
For any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the policy $\hat{\pi}$ returned by VI-LCB using a Bernstein-style penalty term achieves

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

with high prob., with sample complexity at most

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

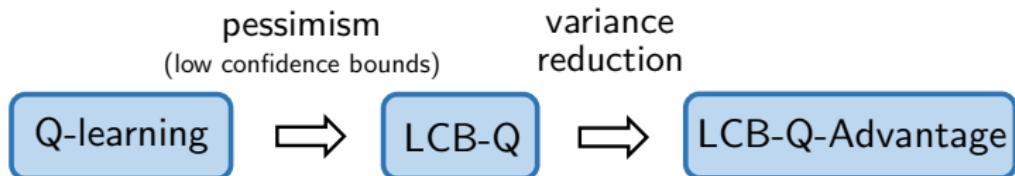
- depends on distribution shift (as reflected by C^*)
- achieves minimax optimality
- full ε -range (no burn-in cost)



Model-based offline RL is minimax optimal with no burn-in cost!

*Is it possible to design offline model-free algorithms
with optimal sample efficiency?*

*Is it possible to design offline model-free algorithms
with optimal sample efficiency?*



LCB-Q: Q-learning with LCB penalty

— Shi et al, 2022, Yan et al, 2023

$$Q_{t+1}(s_t, a_t) \leftarrow \underbrace{(1 - \eta_t)Q_t(s_t, a_t) + \eta_t \mathcal{T}_t(Q_t)(s_t, a_t)}_{\text{classical Q-learning}} - \eta_t \underbrace{b_t(s_t, a_t)}_{\text{LCB penalty}}$$

LCB-Q: Q-learning with LCB penalty

— Shi et al, 2022, Yan et al, 2023

$$Q_{t+1}(s_t, a_t) \leftarrow \underbrace{(1 - \eta_t)Q_t(s_t, a_t) + \eta_t \mathcal{T}_t(Q_t)(s_t, a_t)}_{\text{classical Q-learning}} - \eta_t \underbrace{b_t(s_t, a_t)}_{\text{LCB penalty}}$$

- $b_t(s, a)$: Hoeffding-style confidence bound
- pessimism in the face of uncertainty

LCB-Q: Q-learning with LCB penalty

— Shi et al, 2022, Yan et al, 2023

$$Q_{t+1}(s_t, a_t) \leftarrow \underbrace{(1 - \eta_t)Q_t(s_t, a_t) + \eta_t \mathcal{T}_t(Q_t)(s_t, a_t)}_{\text{classical Q-learning}} - \eta_t \underbrace{b_t(s_t, a_t)}_{\text{LCB penalty}}$$

- $b_t(s, a)$: Hoeffding-style confidence bound
- pessimism in the face of uncertainty

sample size: $\tilde{O}\left(\frac{SC^*}{(1-\gamma)^5 \varepsilon^2}\right) \implies$ sub-optimal by a factor of $\frac{1}{(1-\gamma)^2}$

Issue: large variability in stochastic update rules

Further improvement

pessimism
(low confidence bounds)

Q-learning

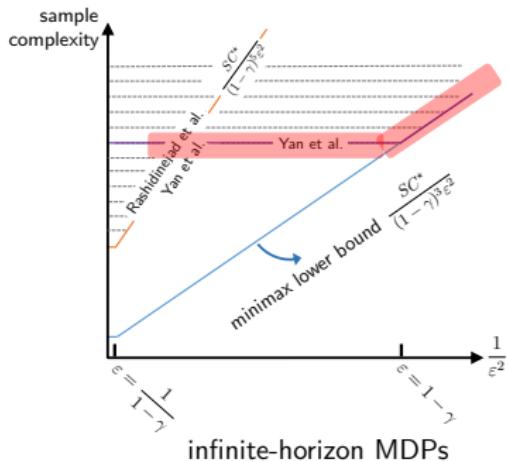


variance reduction

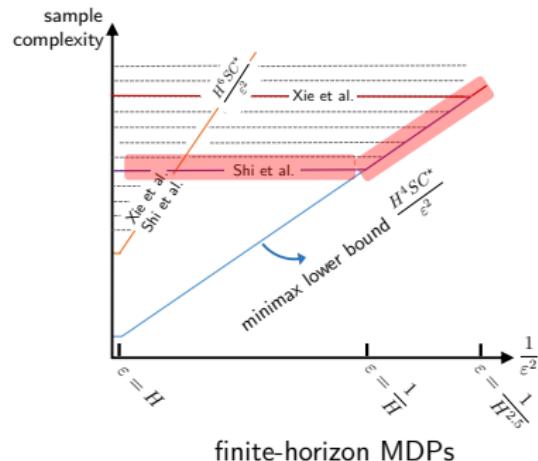
LCB-Q



LCB-Q-Advantage



infinite-horizon MDPs



finite-horizon MDPs

Model-free offline RL attains sample optimality too!

— with some burn-in cost though . . .

Reference: general RL textbooks I

- “*Reinforcement learning: An introduction*,” R. S. Sutton, A. G. Barto, MIT Press, 2018
- “*Reinforcement learning: Theory and algorithms*,” A. Agarwal, N. Jiang, S. Kakade, W. Sun, 2019
- “*Reinforcement learning and optimal control*,” D. Bertsekas, Athena Scientific, 2019
- “*Algorithms for reinforcement learning*,” C. Szepesvari, Springer, 2022
- “*Bandit algorithms*,” T. Lattimore, C. Szepesvari, Cambridge University Press, 2020

Reference: model-based algorithms I

- “*Finite-sample convergence rates for Q-learning and indirect algorithms,*” M. Kearns, S. Satinder, *NeurIPS*, 1998
- “*On the sample complexity of reinforcement learning,*” S. Kakade, 2003
- “*A sparse sampling algorithm for near-optimal planning in large Markov decision processes,*” M. Kearns, Y. Mansour, A. Y. Ng, *Machine learning*, 2002
- “*Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model,*” M. G. Azar, R. Munos, H. J. Kappen, *Machine learning*, 2013
- “*Randomized linear programming solves the Markov decision problem in nearly linear (sometimes sublinear) time,*” *Mathematics of Operations Research*, 2020
- “*Near-optimal time and sample complexities for solving Markov decision processes with a generative model,*” A. Sidford, M. Wang, X. Wu, L. Yang, Y. Ye, *NeurIPS*, 2018

Reference: model-based algorithms II

- “*Variance reduced value iteration and faster algorithms for solving Markov decision processes,*” A. Sidford, M. Wang, X. Wu, Y. Ye, *SODA*, 2018
- “*Model-based reinforcement learning with a generative model is minimax optimal,*” A. Agarwal, S. Kakade, L. Yang, *COLT*, 2020
- “*Instance-dependent ℓ_∞ -bounds for policy evaluation in tabular reinforcement learning,*” A. Pananjady, M. J. Wainwright, *IEEE Trans. on Information Theory*, 2020
- “*Spectral methods for data science: A statistical perspective,*” Y. Chen, Y. Chi, J. Fan, C. Ma, *Foundations and Trends® in Machine Learning*, 2021
- “*Breaking the sample size barrier in model-based reinforcement learning with a generative model,*” G. Li, Y. Wei, Y. Chi, Y. Chen, *Operations Research*, 2024

Reference: model-free algorithms I

- "*A stochastic approximation method,*" H. Robbins, S. Monro, *Annals of Mathematical Statistics*, 1951
- "*Robust stochastic approximation approach to stochastic programming,*" A. Nemirovski, A. Juditsky, G. Lan, A. Shapiro, *SIAM Journal on optimization*, 2009
- "*Q-learning,*" C. Watkins, P. Dayan, *Machine Learning*, 1992
- "*Learning rates for Q-learning,*" E. Even-Dar, Y. Mansour, *Journal of Machine Learning Research*, 2003
- "*Stochastic approximation with cone-contractive operators: Sharp ℓ_∞ bounds for Q-learning,*" M. Wainwright, 2019
- "*Is Q-learning minimax optimal? a tight sample complexity analysis,*" G. Li, C. Cai, Y. Chen, Y. Wei, Y. Chi, *Operations Research*, 2024
- "*Accelerating stochastic gradient descent using predictive variance reduction,*" R. Johnson, T. Zhang, *NeurIPS*, 2013
- "*Variance-reduced Q-learning is minimax optimal,*" M. Wainwright, 2019

Reference: model-free algorithms II

- “*Sample-optimal parametric Q-learning using linearly additive features,*” L. Yang, M. Wang, *ICML*, 2019
- “*Asynchronous stochastic approximation and Q-learning,*” J. Tsitsiklis, *Machine learning*, 1994
- “*Finite-time analysis of asynchronous stochastic approximation and Q-learning,*” G. Qu, A. Wierman, *COLT*, 2020
- “*Finite-sample analysis of contractive stochastic approximation using smooth convex envelopes,*” Z. Chen, S. T. Maguluri, S. Shakkottai, K. Shanmugam, *NeurIPS*, 2020
- “*Sample complexity of asynchronous Q-learning: Sharper analysis and variance reduction,*” G. Li, Y. Wei, Y. Chi, Y. Gu, Y. Chen, *IEEE Trans. on Information Theory*, 2022

Reference: online RL I

- “*Asymptotically efficient adaptive allocation rules,*” T. L. Lai, H. Robbins, *Advances in applied mathematics*, vol. 6, no. 1, 1985
- “*Finite-time analysis of the multiarmed bandit problem,*” P. Auer, N. Cesa-Bianchi, P. Fischer, *Machine learning*, vol. 47, pp. 235-256, 2002
- “*Minimax regret bounds for reinforcement learning,*” M. G. Azar, I. Osband, R. Munos, *ICML*, 2017
- “*Is Q-learning provably efficient?*” C. Jin, Z. Allen-Zhu, S. Bubeck, and M. Jordan, *NeurIPS*, 2018
- “*Provably efficient Q-learning with low switching cost,*” Y. Bai, T. Xie, N. Jiang, Y. X. Wang, *NeurIPS*, 2019
- “*Episodic reinforcement learning in finite MDPs: Minimax lower bounds revisited*” O. D. Domingues, P. Menard, E. Kaufmann, M. Valko, *Algorithmic Learning Theory*, 2021
- “*Almost optimal model-free reinforcement learning via reference-advantage decomposition,*” Z. Zhang, Y. Zhou, X. Ji, *NeurIPS*, 2020

Reference: online RL II

- “*Is reinforcement learning more difficult than bandits? a near-optimal algorithm escaping the curse of horizon,*” Z. Zhang, X. Ji, and S. Du, *COLT*, 2021
- “*Breaking the sample complexity barrier to regret-optimal model-free reinforcement learning,*” G. Li, L. Shi, Y. Chen, Y. Gu, Y. Chi, *NeurIPS*, 2021
- “*Regret-optimal model-free reinforcement learning for discounted MDPs with short burn-in time,*” X. Ji, G. Li, *NeurIPS*, 2023
- “*Reward-free exploration for reinforcement learning,*” C. Jin, A. Krishnamurthy, M. Simchowitz, T. Yu, *ICML*, 2020
- “*Minimax-optimal reward-agnostic exploration in reinforcement learning,*” G. Li, Y. Yan, Y. Chen, J. Fan, *COLT*, 2024
- “*Settling the sample complexity of online reinforcement learning,*” Z. Zhang, Y. Chen, J. D. Lee, S. S. Du, *COLT*, 2024

Reference: offline RL I

- “*Bridging offline reinforcement learning and imitation learning: A tale of pessimism,*” P. Rashidinejad, B. Zhu, C. Ma, J. Jiao, S. Russell, *NeurIPS*, 2021
- “*Is pessimism provably efficient for offline RL?*” Y. Jin, Z. Yang, Z. Wang, *ICML*, 2021
- “*Settling the sample complexity of model-based offline reinforcement learning,*” G. Li, L. Shi, Y. Chen, Y. Chi, Y. Wei, *Annals of Statistics*, vol. 52, no. 1, pp. 233-260, 2024
- “*Pessimistic Q-learning for offline reinforcement learning: Towards optimal sample complexity,*” L. Shi, G. Li, Y. Wei, Y. Chen, Y. Chi, *ICML*, 2022
- “*The efficacy of pessimism in asynchronous Q-learning,*” Y. Yan, G. Li, Y. Chen, J. Fan, *IEEE Transactions on Information Theory*, 2023
- “*Policy finetuning: Bridging sample-efficient offline and online reinforcement learning*” T. Xie, N. Jiang, H. Wang, C. Xiong, Y. Bai, *NeurIPS*, 2021

Recent Advances of Statistical Reinforcement Learning

Part 2

Gergely Neu (Universitat Pompeu Fabra)

Sattar Vakili (MediaTek Research)

Tutorial, UAI 2024

Part 2

- 1. Introduction to structural complexity**
- 2. Linear Function Approximation**
- 3. Non-linear Function Approximation**

Tabular Setting

Recall results for the tabular setting:

- Q-learning with UCB: [Jin et al., 2018]

$$\text{Regret}(T) = \mathcal{O}(\sqrt{H^3SAT})$$

- Sample complexity:

$$\tilde{\mathcal{O}}\left(\frac{\text{poly}(H)SA}{\epsilon^2}\right)$$

Tabular Setting

Recall results for the tabular setting:

- Q-learning with UCB: [Jin et al., 2018]

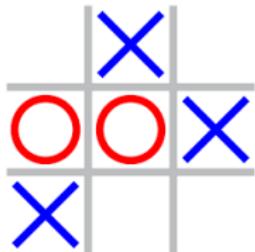
$$\text{Regret}(T) = \mathcal{O}(\sqrt{H^3SAT})$$

- Sample complexity:

$$\tilde{\mathcal{O}}\left(\frac{\text{poly}(H)SA}{\epsilon^2}\right)$$

These are only meaningful if $T \ll S$ or $\epsilon \gg 1/\sqrt{S}$!

Why “Tabular”?



- Small size of state-action space
- $Q(s, a)$ can be represented as a [table](#)

Why Function Approximation?

Number of states S is **enormous** in real-world problems!



- Game of Go: 10^{170} states
- Atari: 10^{100} states
- Physical systems:
continuum of states

Why Function Approximation?



Two types of challenges:

- ▶ **Computational:** Q and π cannot even be stored in memory, and Bellman equations are intractable to solve even if P and r were known
- ▶ **Statistical:** Most states are not visited even once! How could we expect to learn about P or r like that?

Why Function Approximation?



Two types of challenges:

- ▶ **Computational:** Q and π cannot even be stored in memory, and Bellman equations are intractable to solve even if P and r were known
- ▶ **Statistical:** Most states are not visited even once! How could we expect to learn about P or r like that?

We need to find a way to **generalize** knowledge from visited states to unvisited states by leveraging **structure**

RL with Function Approximation

- ▶ Approximate value function $Q(s, a)$ (or policy) in a class \mathcal{F} .
- ▶ Hope that \mathcal{F} captures the MDP structure appropriately and leverage the information in structure of \mathcal{F} to learn faster if possible.
- ▶ Typical function classes: Linear, Kernel-based, NN-based

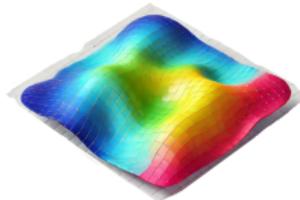
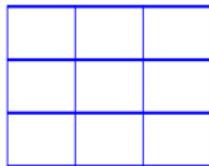
Tabular → Linear → Nonlinear

Setting

- ▶ Generative oracle, Offline, Online
- ▶ Episodic, Infinite horizon (discounted)
- ▶ Model-based, Model-free

In this part we focus on:

Tabular → Linear → Nonlinear

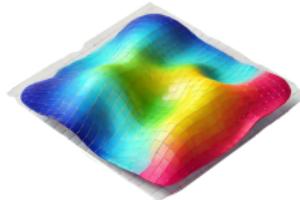
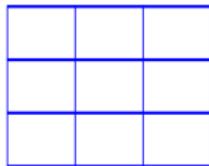


Setting

- ▶ Generative oracle, Offline, Online
- ▶ Episodic, Infinite horizon (discounted)
- ▶ Model-based, Model-free

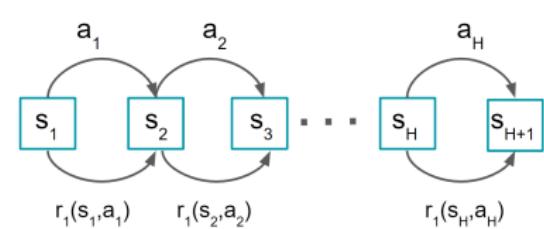
In this part we focus on:

Tabular → Linear → Nonlinear



Setting

For a clear and sharp presentation we focus on **episodic MDPs**

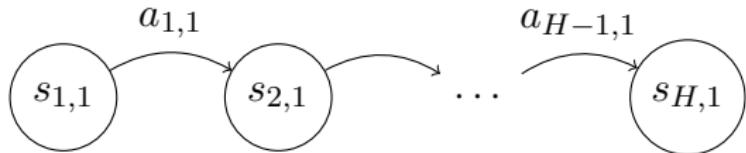


For simplicity, we assume r is known and deterministic

We focus on the structural complexity of $P(s'|s, a)$

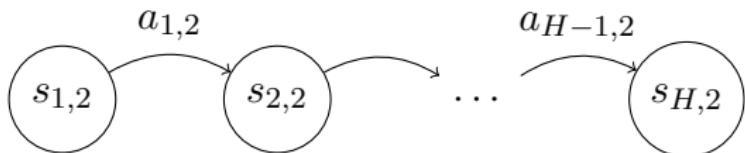
Episodic MDP

Episode 1:

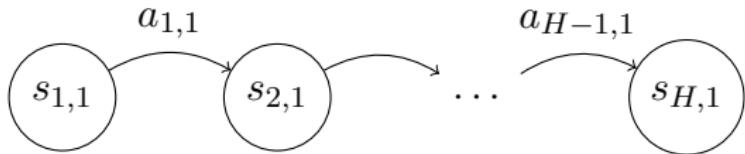


Episodic MDP

Episode 2:

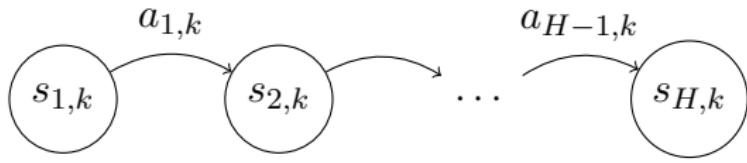


Episode 1:



Episodic MDP

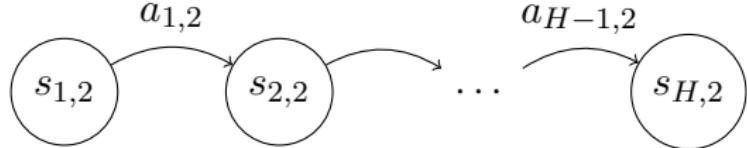
Episode k :



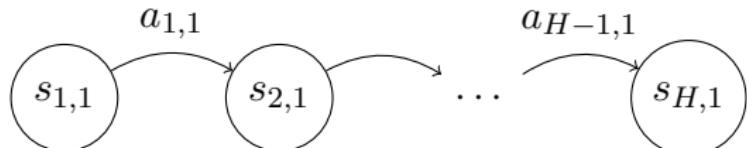
⋮

⋮

Episode 2:



Episode 1:



Part 2

1. Introduction to structural complexity
2. Linear Function Approximation
3. Non-linear Function Approximation

Linear Function Approximation

IDEA: approximate the Q -functions as linear functions of a given d -dimensional feature map $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$.

Linear Function Approximation

IDEA: approximate the Q -functions as linear functions of a given d -dimensional feature map $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$.

Let $\Phi \in \mathbb{R}^{(\mathcal{S} \times \mathcal{A}) \times d}$ be the “matrix” of stacked feature vectors $[\phi(s_1, a_1) \dots \phi(s_N, a_N)]^\top$ (where $N = |\mathcal{S} \times \mathcal{A}|$).

We need to find a parameter vector θ^* such that $Q^* \approx \Phi\theta^*$.
(Meaning that $Q^*(s, a) \approx \langle \theta^*, \phi(s, a) \rangle$.)

Linear Function Approximation

IDEA: approximate the Q -functions as linear functions of a given d -dimensional feature map $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$.

Let $\Phi \in \mathbb{R}^{(\mathcal{S} \times \mathcal{A}) \times d}$ be the “matrix” of stacked feature vectors $[\phi(s_1, a_1) \dots \phi(s_N, a_N)]^\top$ (where $N = |\mathcal{S} \times \mathcal{A}|$).

We need to find a parameter vector θ^* such that $Q^* \approx \Phi\theta^*$.
(Meaning that $Q^*(s, a) \approx \langle \theta^*, \phi(s, a) \rangle$.)

QUESTION: When can we learn θ^* efficiently?

Linear Function Approximation

Various conditions on the feature map Φ have been studied:

- **Linear Q^* :** there exists a θ^* such that $Q^* = \Phi\theta^*$.
- **Linear Q^π :** for every policy π , there exists a θ^π such that $Q^\pi = \Phi\theta^\pi$.
- **Closure under Bellman operator:** for any $Q_\theta = \Phi\theta$,
 $\mathcal{T}Q_\theta \in \text{span}(\Phi)$.
- **Linear MDP:** The transition and reward functions are linear in
the features. This implies all of the above conditions.

A number of more refined conditions have been also studied, such as assuming linearity of V^* in some feature map, or other types of factorized transition models. We refer to [Du et al. \[2021\]](#), [Jin et al. \[2021\]](#) for more details on such extensions.

What can we hope for?

WANT: find an ϵ -optimal policy with a computational and sample complexity polynomial in d , $1/\epsilon$ and H , independently of S and A .

What can we hope for?

WANT: find an ϵ -optimal policy with a computational and sample complexity polynomial in d , $1/\epsilon$ and H , independently of S and A .

- This is **impossible** when only requiring linear Q^* -realizability!
[Weisz et al., 2021]

What can we hope for?

WANT: find an ϵ -optimal policy with a computational and sample complexity polynomial in d , $1/\epsilon$ and H , independently of S and A .

- This is **impossible** when only requiring linear Q^* -realizability! [Weisz et al., 2021]
- Polynomial sample complexity is possible when relaxing the condition to linear Q^π -realizability, but no practical algorithms are known [Weisz et al., 2022, 2023].

What can we hope for?

WANT: find an ϵ -optimal policy with a computational and sample complexity polynomial in d , $1/\epsilon$ and H , independently of S and A .

- This is **impossible** when only requiring linear Q^* -realizability! [Weisz et al., 2021]
- Polynomial sample complexity is possible when relaxing the condition to linear Q^π -realizability, but no practical algorithms are known [Weisz et al., 2022, 2023].
- Situation is similar when only assuming closure under Bellman operator / Bellman completeness [Zanette et al., 2020b, Du et al., 2021].

What can we hope for?

WANT: find an ϵ -optimal policy with a computational and sample complexity polynomial in d , $1/\epsilon$ and H , independently of S and A .

- This is **impossible** when only requiring linear Q^* -realizability! [Weisz et al., 2021]
- Polynomial sample complexity is possible when relaxing the condition to linear Q^π -realizability, but no practical algorithms are known [Weisz et al., 2022, 2023].
- Situation is similar when only assuming closure under Bellman operator / Bellman completeness [Zanette et al., 2020b, Du et al., 2021].
- Linear MDP condition enables both statistical and computational efficiency!!! [Jin et al., 2023]

Linear MDPs

Linear transition function:

$$P_h(\cdot|s, a) = \langle \phi(s, a), \mu_h(\cdot) \rangle,$$

where $\mu_h(\cdot) = [\mu_h^1(\cdot), \dots, \mu_h^d(\cdot)]$ is a d -dimensional signed measure.

Linear rewards: $r_h(s, a) = \langle \phi(s, a), \vartheta_h \rangle$.

Linear MDPs

Linear transition function:

$$P_h(\cdot|s, a) = \langle \phi(s, a), \mu_h(\cdot) \rangle,$$

where $\mu_h(\cdot) = [\mu_h^1(\cdot), \dots, \mu_h^d(\cdot)]$ is a d -dimensional signed measure.

Linear rewards: $r_h(s, a) = \langle \phi(s, a), \vartheta_h \rangle$.

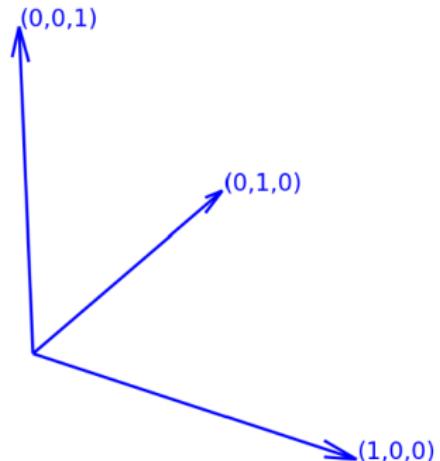
In matrix notation:

- Transition operator $P_h \in \mathbb{R}^{(\mathcal{S} \times \mathcal{A}) \times \mathcal{S}}$ can be written as $P_h = \Phi M_h$ for some “matrix” $M_h \in \mathbb{R}^{\mathcal{S} \times d}$.
- Reward function can be written as $r_h = \Phi \vartheta_h$ for some $\vartheta_h \in \mathbb{R}^d$.

Tabular MDPs are linear

Tabular setting is a special case with dimension $d = SA$:

- Let $\phi(s, a) = e_{(s,a)}$ be the canonical basis in \mathbb{R}^d
- $P_h(\cdot|s, a) = e_{s,a}^\top \mu_h(\cdot)$



A magical property of linear MDPs

In a linear MDP, the Q -functions of all policies are linear in Φ :

$$\begin{aligned} Q_h^\pi &= r_h + P_h V_{h+1}^\pi = \Phi \vartheta_h + \Phi M_h V_{h+1}^\pi \\ &= \Phi (\vartheta_h + M_h V_{h+1}^\pi) = \Phi \boldsymbol{\theta}_h, \end{aligned}$$

with $\boldsymbol{\theta}_h = \vartheta_h + M_h V_{h+1}^\pi$.

A magical property of linear MDPs

In a linear MDP, the Q -functions of all policies are linear in Φ :

$$\begin{aligned} Q_h^\pi &= r_h + P_h V_{h+1}^\pi = \Phi \vartheta_h + \Phi M_h V_{h+1}^\pi \\ &= \Phi (\vartheta_h + M_h V_{h+1}^\pi) = \Phi \boldsymbol{\theta}_h, \end{aligned}$$

with $\boldsymbol{\theta}_h = \vartheta_h + M_h V_{h+1}^\pi$.

This implies linear Q^* -realizability, linear Q^π -realizability, Bellman completeness, and many more useful properties for analysis! E.g., note that for any function $u \in \mathbb{R}^S$, $P_h u = \Phi M_h u$ is linear in Φ .

A magical property of linear MDPs

In a linear MDP, the Q -functions of all policies are linear in Φ :

$$\begin{aligned} Q_h^\pi &= r_h + P_h V_{h+1}^\pi = \Phi \vartheta_h + \Phi M_h V_{h+1}^\pi \\ &= \Phi (\vartheta_h + M_h V_{h+1}^\pi) = \Phi \boldsymbol{\theta}_h, \end{aligned}$$

with $\boldsymbol{\theta}_h = \vartheta_h + M_h V_{h+1}^\pi$.

This implies linear Q^* -realizability, linear Q^π -realizability, Bellman completeness, and many more useful properties for analysis! E.g., note that for any function $u \in \mathbb{R}^S$, $P_h u = \Phi M_h u$ is linear in Φ .

The structure of linear MDPs allows us to import tools from linear bandit literature [Abbasi-Yadkori et al., 2011, Lattimore and Szepesvári, 2020].

Optimistic approximate dynamic programming

IDEA: Combine the techniques for tabular MDPs with exploration bonuses borrowed from the linear bandit literature!

Optimistic approximate dynamic programming

IDEA: Combine the techniques for tabular MDPs with exploration bonuses borrowed from the linear bandit literature!

UCB-VI [Azar et al., 2017]:

- ① Backtrack $h = H, H - 1, \dots, 1$: run optimistic value iteration

$$Q_h = r_h + \underbrace{\hat{P}_h}_{\text{model estimate}} V_{h+1} + \underbrace{b_h}_{\text{exploration bonus}}$$

and set $V_h(s) = \max_a Q_h(s, a)$ for all s, a .

- ② Forward $h = 1, 2, \dots, H$: take actions according to greedy policy

$$\pi_h(s) = \arg \max_a Q_h(s, a).$$

Optimistic approximate dynamic programming

IDEA: Combine the techniques for tabular MDPs with exploration bonuses borrowed from the linear bandit literature!

UCB-VI [Azar et al., 2017]:

- ① Backtrack $h = H, H - 1, \dots, 1$: run optimistic value iteration

$$Q_h = r_h + \underbrace{\hat{P}_h}_{\text{model estimate}} V_{h+1} + \underbrace{b_h}_{\text{exploration bonus}}$$

and set $V_h(s) = \max_a Q_h(s, a)$ for all s, a .

- ② Forward $h = 1, 2, \dots, H$: take actions according to greedy policy

$$\pi_h(s) = \arg \max_a Q_h(s, a).$$

But how do we define \hat{P}_h and b_h ?

Least Squares Value Iteration (LSVI)

Transition model \widehat{P}_h can be defined **implicitly** via least-squares:

- ▶ Solve the regularized linear regression problem

$$\widehat{\mathbf{w}}_{h,k} = \arg \min_{\mathbf{w}} \sum_{t=1}^k (V_{h+1,k}(s_{h+1,t}) - \langle \phi(s_{h,t}, a_{h,t}), \mathbf{w} \rangle)^2 + \lambda^2 \|\mathbf{w}\|^2$$

- ▶ That provides a prediction

$$[\widehat{P}_h V_{h+1,k}](s, a) = \langle \phi(s, a), \widehat{\mathbf{w}}_{h,k} \rangle$$

- ▶ Also, an uncertainty quantification (variance)

$$\sigma_{h,k}^2(s, a) = \|\phi(s, a)\|_{(\lambda I + \Sigma_{h,k})^{-1}}^2$$

$$\Sigma_{h,k} = \sum_{t=1}^k \phi^\top(s_{h,t}, a_{h,t}) \phi(s_{h,t}, a_{h,t})$$

LSVI-UCB

The prediction and variance give us an upper confidence bound on Q^* :

$$Q_{h,k}(s, a) = r_h(s, a) + \widehat{[P_h V_{h+1}]}(s, a) + \beta(\delta) \sigma_{h,k}(s, a)$$

This is then used to compute an UCB on V^* as

$$V_{h,k}(s) = \max_a Q_{h,k}(s, a)$$

Performance guarantees

Theorem [Jin et al., 2023] The regret of LSVI-UCB satisfies
 $\text{Regret}(K) = \tilde{\mathcal{O}}(H^2\sqrt{d^3K}).$

This implies a sample complexity guarantee of $\tilde{\mathcal{O}}(\frac{\text{poly}(H)d^3}{\epsilon^2})$.

Performance guarantees

Theorem [Jin et al., 2023] The regret of LSVI-UCB satisfies
 $\text{Regret}(K) = \tilde{\mathcal{O}}(H^2\sqrt{d^3 K})$.

This implies a sample complexity guarantee of $\tilde{\mathcal{O}}(\frac{\text{poly}(H)d^3}{\epsilon^2})$.

Proof ideas:

- Prove confidence bounds

$$|\widehat{[P_h V_{h+1,k}]}(s, a) - [P_h V_{h+1,k}](s, a)| \leq \beta(\delta) \sigma_{h,k}(s, a).$$

- Using standard techniques (e.g., Azar et al., 2017), show

$$\text{Regret}(K) \lesssim \sum_{h,k} \beta(\delta) \sigma_{h,k}(s_{h,k}, a_{h,k}).$$

- Use elliptical potential lemma (e.g., Abbasi-Yadkori et al., 2011) to show

$$\sum_{h,k} \sigma_{h,k}(s_{h,k}, a_{h,k}) \lesssim H \sqrt{K d \log(K)}.$$

Proving the confidence bounds

By standard results on least-squares estimators (e.g., Abbasi-Yadkori et al., 2011), one can prove the following confidence bound for any **fixed** $u \in \mathbb{R}^S$ that holds with probability at least $1 - \delta$:

$$\left| \widehat{[P_h u]}(s, a) - [P_h u](s, a) \right| \leq \bar{\beta}(\delta) \sigma_{h,k}(s, a)$$

for some

$$\bar{\beta}(\delta) \approx \lambda^{\frac{1}{2}} \|M_h u\| + H \sqrt{d \log\left(\frac{K}{\delta}\right)}$$

Proving the confidence bounds

By standard results on least-squares estimators (e.g., Abbasi-Yadkori et al., 2011), one can prove the following confidence bound for any **fixed** $u \in \mathbb{R}^S$ that holds with probability at least $1 - \delta$:

$$\left| \widehat{[P_h u]}(s, a) - [P_h u](s, a) \right| \leq \overline{\beta}(\delta) \sigma_{h,k}(s, a)$$

for some

$$\overline{\beta}(\delta) \approx \lambda^{\frac{1}{2}} \|M_h u\| + H \sqrt{d \log\left(\frac{K}{\delta}\right)}$$

Challenge: $u = V_{h+1,k}$ is **not** fixed, but depends on all past data!

Proving the confidence bounds

By standard results on least-squares estimators (e.g., Abbasi-Yadkori et al., 2011), one can prove the following confidence bound for any **fixed** $u \in \mathbb{R}^S$ that holds with probability at least $1 - \delta$:

$$\left| \widehat{[P_h u]}(s, a) - [P_h u](s, a) \right| \leq \bar{\beta}(\delta) \sigma_{h,k}(s, a)$$

for some

$$\bar{\beta}(\delta) \approx \lambda^{\frac{1}{2}} \|M_h u\| + H \sqrt{d \log\left(\frac{K}{\delta}\right)}$$

Challenge: $u = V_{h+1,k}$ is **not** fixed, but depends on all past data!

Solution: Covering number argument

Covering Number Argument

- ▶ Notice that all value functions $V_{h,k}$ belong to the function class

$$\mathcal{V} = \left\{ V(s) = \min\{H, \max_a \phi^\top(s, a)\hat{w} + \beta\|\phi(s, a)\|_{(\lambda I + \Sigma)^{-1}}\} \right\}.$$

- ▶ **Idea:** cover the space of functions \mathcal{V} such that we can rewrite

$$\left| \widehat{[P_h V_{h+1,k}]}(s, a) - [P_h V_{h+1,k}](s, a) \right| \leq \epsilon + \sup_{u \in \mathcal{V}} \left| \widehat{[P_h u]}(s, a) - [P_h u](s, a) \right|.$$

- ▶ How many functions u are required to cover \mathcal{V} up to ϵ error?

$$\mathcal{N}_\epsilon = \tilde{\mathcal{O}}(d^2)$$

Covering Number Argument

- ▶ Notice that all value functions $V_{h,k}$ belong to the function class

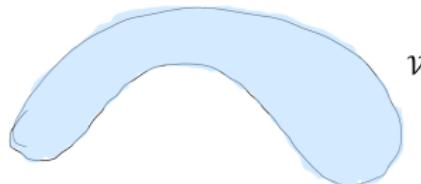
$$\mathcal{V} = \left\{ V(s) = \min\{H, \max_a \phi^\top(s, a)\hat{w} + \beta\|\phi(s, a)\|_{(\lambda I + \Sigma)^{-1}}\} \right\}.$$

- ▶ **Idea:** cover the space of functions \mathcal{V} such that we can rewrite

$$\left| \widehat{[P_h V_{h+1,k}]}(s, a) - [P_h V_{h+1,k}](s, a) \right| \leq \epsilon + \sup_{u \in \mathcal{V}} \left| \widehat{[P_h u]}(s, a) - [P_h u](s, a) \right|.$$

- ▶ How many functions u are required to cover \mathcal{V} up to ϵ error?

$$\mathcal{N}_\epsilon = \tilde{\mathcal{O}}(d^2)$$



Covering Number Argument

- ▶ Notice that all value functions $V_{h,k}$ belong to the function class

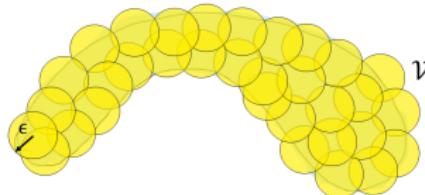
$$\mathcal{V} = \left\{ V(s) = \min\{H, \max_a \phi^\top(s, a)\hat{w} + \beta\|\phi(s, a)\|_{(\lambda I + \Sigma)^{-1}}\} \right\}.$$

- ▶ **Idea:** cover the space of functions \mathcal{V} such that we can rewrite

$$\left| \widehat{[P_h V_{h+1,k}]}(s, a) - [P_h V_{h+1,k}](s, a) \right| \leq \epsilon + \sup_{u \in \mathcal{V}} \left| \widehat{[P_h u]}(s, a) - [P_h u](s, a) \right|.$$

- ▶ How many functions u are required to cover \mathcal{V} up to ϵ error?

$$\mathcal{N}_\epsilon = \tilde{\mathcal{O}}(d^2)$$



Covering Number Argument

- ▶ We can now use a union-bound argument to show that

$$\sup_{u \in \mathcal{V}} \left| \widehat{[P_h u]}(s, a) - [P_h u](s, a) \right| \leq \epsilon + \bar{\beta}(\delta / \mathcal{N}_\epsilon) \sigma_{h,k}(s, a)$$

holds with probability at least $1 - \delta$.

- ▶ Choosing $\epsilon \approx 1/T$, we get

$$\beta(\delta) = \bar{\beta}(\delta / \mathcal{N}_\epsilon) \approx \lambda^{\frac{1}{2}} \|MV_{h,k}\| + H \sqrt{d \log\left(\frac{K \mathcal{N}_\epsilon}{\delta}\right)} = \tilde{\mathcal{O}}(Hd)$$

Alternative linear models

Linear MDP model factorizes $P = \Phi M$ with some known $\Phi \in \mathbb{R}^{(S \times A) \times d}$ and some unknown $M \in \mathbb{R}^{d \times S}$.

Some alternative factorizations are:

- **Linear mixture MDPs** [Zhou et al., 2021]: $P = \Phi \theta$ with some known $\Phi \in \mathbb{R}^{(S \times A \times S) \times d}$ and unknown $\theta \in \mathbb{R}^d$. Analysis is simpler but the model doesn't allow simple and explicit Q -function approximation and leads to impractical algorithms.
- **"MatrixRL"** [Yang and Wang, 2020]: $P = \Phi M \Psi$ with some known $\Phi \in \mathbb{R}^{(S \times A) \times m}$, another known $\Psi \in \mathbb{R}^{n \times S}$, and an unknown $M \in \mathbb{R}^{m \times n}$. Can be shown to be a special case of linear mixture MDPs, and suffers from the same limitations.
- **Low-rank MDPs** [Modi et al., 2024]: Same as linear MDPs except both Φ and M are unknown and belong to finite model class. Requires much more sophisticated techniques, but algorithms are kind of tractable.

Some References

- Linear MDPs: Jin et al. [2020, 2023], Yang and Wang [2019, 2020], Neu and Pike-Burke [2020]
- Linear Bellman complete models: Zanette et al. [2020a]
- Linear mixture MDPs: Yang and Wang [2020], Ayoub et al. [2020], Zhou et al. [2021], Moulin and Neu [2023]
- Other model classes with hidden finite-dimensional linear structure: Du et al. [2021], Jin et al. [2021]

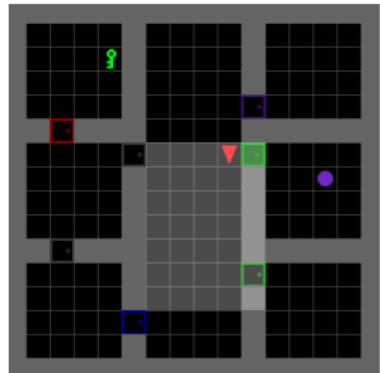
Part 2

- 1.** Introduction to structural complexity
- 2.** Linear Function Approximation
- 3.** Non-linear Function Approximation

Limitations of the Linear Setting

Directly reachable states:

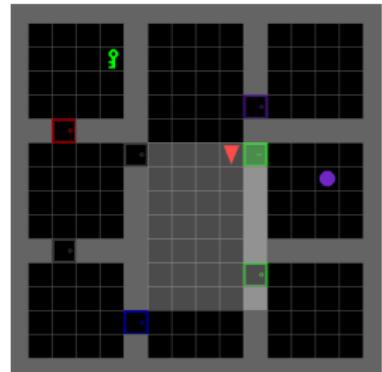
- $\mathcal{S}_{s,a} := \{s' \in \mathcal{S} : P(s'|s, a) > 0\}$
- $U := \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} |\mathcal{S}_{s,a}|$



Limitations of the Linear Setting

Directly reachable states:

- $\mathcal{S}_{s,a} := \{s' \in \mathcal{S} : P(s'|s, a) > 0\}$
- $U := \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} |\mathcal{S}_{s,a}|$



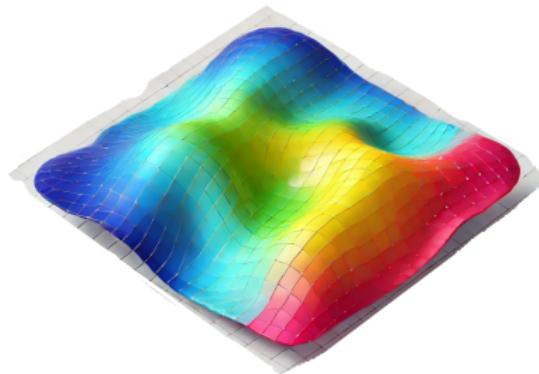
Theorem Lee and Oh [2024] For an MDP with a finite state space, the feature dimension d is lower bounded by

$$d \geq \lfloor \frac{|\mathcal{S}|}{U} \rfloor$$

Limitations of the Linear Setting



High dimensional problems



Nonlinear problems

Kernel-Based Setting

Kernel-Based Setting

- Kernel-based models are natural extensions of linear models to **infinite dimensional** feature maps

Kernel-Based Setting

- Kernel-based models are natural extensions of linear models to **infinite dimensional** feature maps
- Allow for versatile and powerful **non-linear function approximation**

Kernel-Based Setting

- Kernel-based models are natural extensions of linear models to **infinite dimensional** feature maps
- Allow for versatile and powerful **non-linear function approximation**
- Lend themselves to analysis

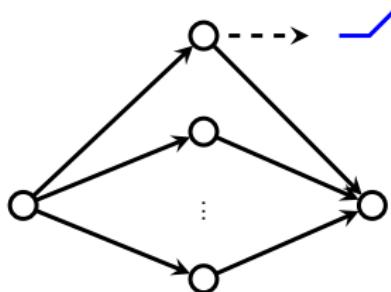
Kernel-Based Setting

- Kernel-based models are natural extensions of linear models to **infinite dimensional** feature maps
- Allow for versatile and powerful **non-linear function approximation**
- Lend themselves to analysis
- Serve as an intermediate step towards analysis of NN-based models

Kernel-Based Setting

- Kernel-based models are natural extensions of linear models to **infinite dimensional** feature maps
- Allow for versatile and powerful **non-linear function approximation**
- Lend themselves to analysis
- Serve as an intermediate step towards analysis of NN-based models

Tabular → Linear → Kernel-Based → NN-Based



Kernel-Based Setting

Function class:

$$\mathcal{F} = \left\{ f : \mathbb{R}^d \rightarrow \mathbb{R}, \ f(\cdot) = \sum_{m=1}^{\infty} w_m \phi_m(\cdot) \right\}$$

Kernel-Based Setting

Function class:

$$\mathcal{F} = \left\{ f : \mathbb{R}^d \rightarrow \mathbb{R}, \ f(\cdot) = \sum_{m=1}^{\infty} w_m \phi_m(\cdot) \right\}$$

An extension of linear models to infinite dimensions in the feature space ϕ

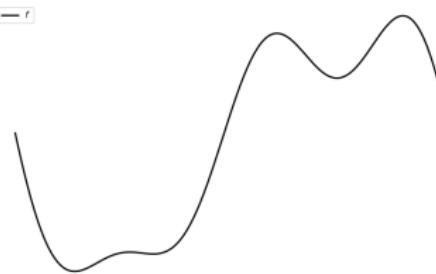
Kernel-Based Setting

Function class:

$$\mathcal{F} = \left\{ f : \mathbb{R}^d \rightarrow \mathbb{R}, f(\cdot) = \sum_{m=1}^{\infty} w_m \phi_m(\cdot) \right\}$$

An extension of linear models to infinite dimensions in the feature space ϕ

Nonlinear functions in \mathbb{R}^d



Mercer Theorem

A positive definite kernel $\kappa : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$

Mercer Theorem

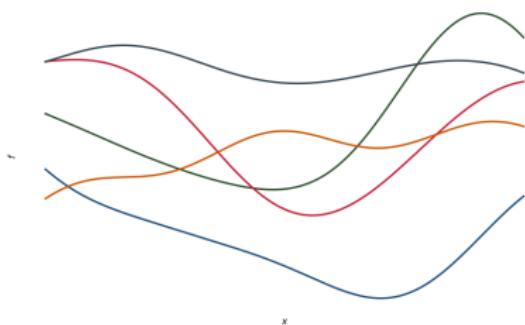
A positive definite kernel $\kappa : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$

Theorem Any positive definite kernel can be written as

$$\kappa(z, z') = \sum_{m=1}^{\infty} \lambda_m \varphi_m(z) \varphi_m(z')$$

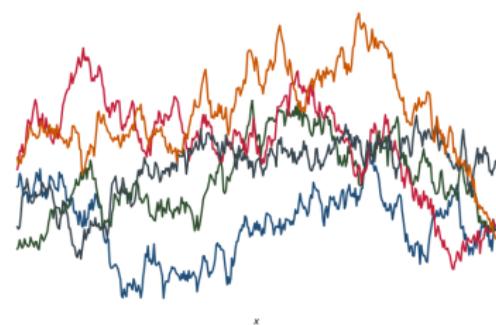
- The feature map $\phi_m(\cdot) = \lambda_m^{\frac{1}{2}} \varphi_m(\cdot)$ corresponding to κ
- λ_m are referred to as eigenvalues
- φ_m are referred to as eigenfunctions

Kernels



Squared Exponential kernel

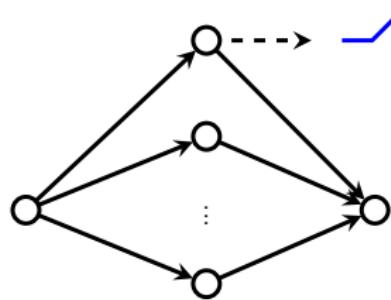
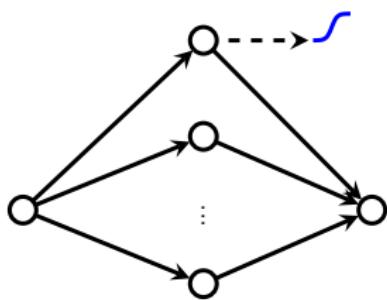
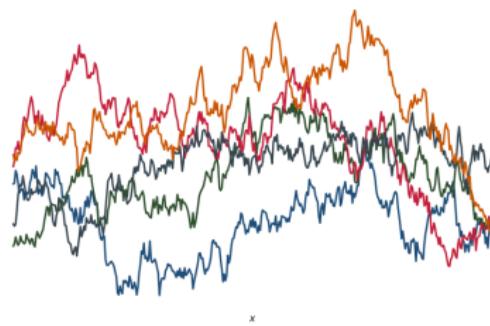
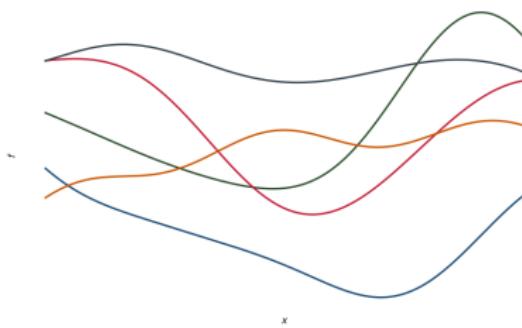
$$\kappa(z, z') = \exp\left(-\frac{\|z-z'\|^2}{2\ell^2}\right)$$



Matérn- ν kernel

$$\kappa(z, z') = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}}{\ell} \|z - z'\|\right)^\nu K_\nu \left(\frac{\sqrt{2\nu}}{\ell} \|z - z'\|\right)$$

Kernels



Reproducing Kernel Hilbert Space

RKHS:

$$\mathcal{H}_\kappa = \{f(\cdot) = \sum_{m=1}^{\infty} w_m \phi_m(\cdot)\}$$

- Inner product $\langle f, g \rangle_k = \mathbf{w}_f^\top \mathbf{w}_g$
- $\|f\|_{\mathcal{H}_\kappa} = \|\mathbf{w}\|$
- $\phi_m = \sqrt{\lambda_m} \varphi_m$ form an orthonormal basis

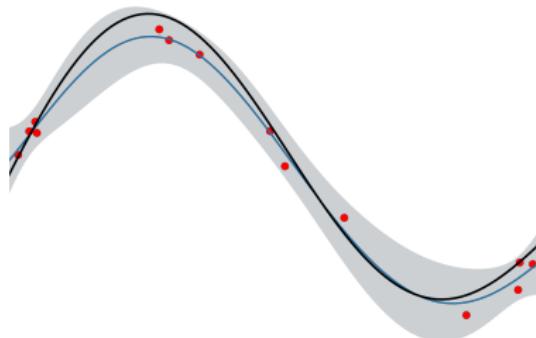
Kernel Based Regression

Provided a dataset of t observation:

$$\left\{ (z_j, Y(z_j)) \right\}_{j=1}^t, Y(z_j) = f(z_j) + \varepsilon_j$$

Regularized Least Squares Error:

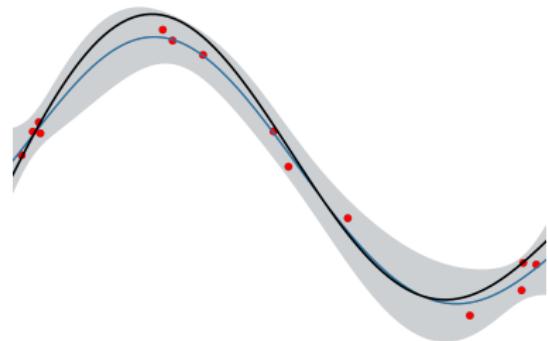
$$\hat{f} = \arg \min_{g \in \mathcal{H}_\kappa} \sum_{j=1}^t (Y(z_j) - g(z_j))^2 + \lambda \|g\|_{\mathcal{H}_\kappa}^2$$



Kernel-Based Regression

Predictor:

$$\hat{f}(z) = \kappa_t^\top(z)(\mathbf{K}_t + \lambda I)^{-1}\mathbf{y}_t$$

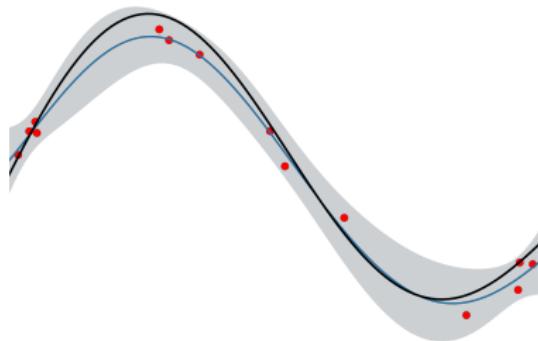


- $\kappa_t(z) = [k(z_1, z), k(z_2, z), \dots, k(z_t, z)]$
- $\mathbf{K}_t = [k(z_i, z_j)]_{i,j=1}^t$
- $\mathbf{y}_t = [Y(z_1), Y(z_2), \dots, Y(z_t)]$

Kernel-Based Regression

Uncertainty estimator:

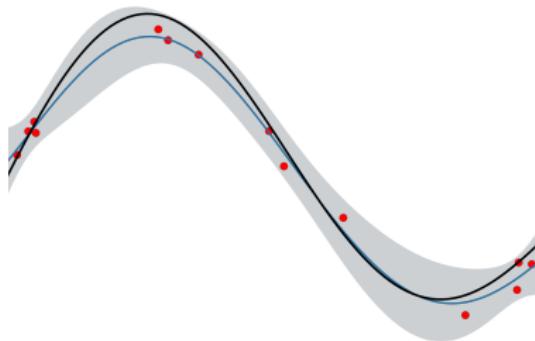
$$(\sigma_t(z))^2 = \kappa(z, z) - \boldsymbol{\kappa}_t^\top(z)(\mathbf{K}_t + \lambda I)^{-1}\boldsymbol{\kappa}_t(z)$$



Kernel-Based Regression

Uncertainty estimator:

$$(\sigma_t(z))^2 = \kappa(z, z) - \boldsymbol{\kappa}_t^\top(z)(\mathbf{K}_t + \lambda I)^{-1}\boldsymbol{\kappa}_t(z)$$



Closed form expressions for prediction and uncertainty quantification!

RL with kernel-based function approximation

IDEA: Approximate the Q -functions as a function in RKHS

RL with kernel-based function approximation

IDEA: Approximate the Q -functions as a function in RKHS

Want: Find an ϵ -optimal policy with a computational and sample complexity polynomial in $1/\epsilon$ and H

RL with kernel-based function approximation

IDEA: Approximate the Q -functions as a function in RKHS

Want: Find an ϵ -optimal policy with a computational and sample complexity polynomial in $1/\epsilon$ and H

—possibly some kernel parameters—

RL with kernel-based function approximation

IDEA: Approximate the Q -functions as a function in RKHS

Want: Find an ϵ -optimal policy with a computational and sample complexity polynomial in $1/\epsilon$ and H

—possibly some kernel parameters—
independently of S and A .

RL with kernel-based function approximation

IDEA: Approximate the Q -functions as a function in RKHS

Want: Find an ϵ -optimal policy with a computational and sample complexity polynomial in $1/\epsilon$ and H

—possibly some kernel parameters—
independently of S and A .

Sample complexity: $\tilde{\mathcal{O}}\left((\frac{1}{\epsilon})^2\right)$

RL with kernel-based function approximation

IDEA: Approximate the Q -functions as a function in RKHS

Want: Find an ϵ -optimal policy with a computational and sample complexity polynomial in $1/\epsilon$ and H

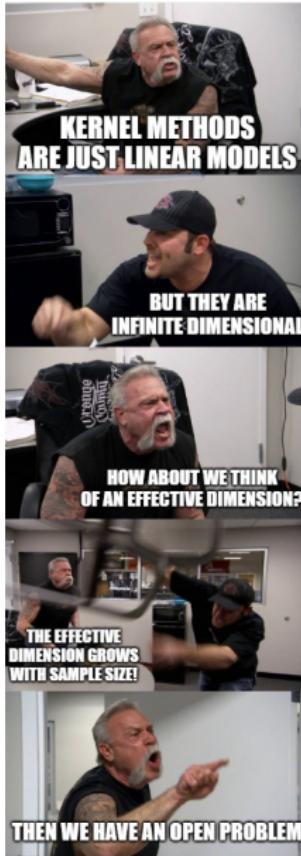
—possibly some kernel parameters—
independently of S and A .

Sample complexity: $\tilde{\mathcal{O}}\left((\frac{1}{\epsilon})^?\right)$

RL with kernel-based function approximation



RL with kernel-based function approximation



Effective Dimension:

$$\left[\underbrace{\phi_1, \phi_2, \dots, \phi_{\mathfrak{D}}}_{\mathfrak{D}\text{dimension}}, \phi_{\mathfrak{D}+1}, \dots \right]$$

$$\mathfrak{D} \approx \frac{1}{2} \log \det(I + \frac{1}{\lambda} \mathbf{K}_t)$$

- In the linear setting: $\mathfrak{D} \approx d$
- For Squared Exponential kernel:
 $\mathfrak{D} \approx \text{poly log}(T)$
- For Matérn kernel:
 $\mathfrak{D} \approx T^{\frac{d}{d+\nu}}$ [Vakili et al., 2021]

RL with Kernel-based FA

Kernel-based transition assumption

For all s' : $P(s'|s, a) \in \mathcal{H}_\kappa$

RL with Kernel-based FA

Kernel-based transition assumption

For all s' : $P(s'|s, a) \in \mathcal{H}_\kappa$

- A significant generalization of linear models

RL with Kernel-based FA

Kernel-based transition assumption

For all s' : $P(s'|s, a) \in \mathcal{H}_\kappa$

- A significant generalization of linear models
- Linear model is a special case with linear kernel:
 $\kappa(s, a, s', a') = \phi^\top(s, a)\phi(s', a')$

RL with Kernel-based FA

Kernel-based transition assumption

For all s' : $P(s'|s, a) \in \mathcal{H}_\kappa$

- A significant generalization of linear models
- Linear model is a special case with linear kernel:
 $\kappa(s, a, s', a') = \phi^\top(s, a)\phi(s', a')$
- RKHS of common kernels can approximate almost all continuous functions

RL with Kernel-based FA

Kernel-based transition assumption

For all s' : $P(s'|s, a) \in \mathcal{H}_\kappa$

- A significant generalization of linear models
- Linear model is a special case with linear kernel:
 $\kappa(s, a, s', a') = \phi^\top(s, a)\phi(s', a')$
- RKHS of common kernels can approximate almost all continuous functions

For integrable $V : \mathcal{S} \rightarrow \mathbb{R}$, $[PV] = \int_{s'} P(s'|s, a)V(s') \in \mathcal{H}_k$

Optimistic approximate DP goes kernelized

IDEA: Combine the techniques for tabular MDPs with exploration bonuses borrowed from the linear kernel bandit literature!

UCB-VI [Azar et al., 2017]:

- ① Backtrack $h = H, H - 1, \dots, 1$: run optimistic value iteration

$$Q_h = r_h + \underbrace{\hat{P}_h}_{\text{model estimate}} V_{h+1} + \underbrace{b_h}_{\text{exploration bonus}}$$

and set $V_h(s) = \max_a Q_h(s, a)$ for all s, a .

- ② Forward $h = 1, 2, \dots, H$: take actions according to greedy policy

$$\pi_h(s) = \arg \max_a Q_h(s, a).$$

But how do we define \hat{P}_h and b_h ?

Kernel-based optimistic value iteration (KOVI)

Transition model \widehat{P}_h can be defined **implicitly** via least-squares:

- ▶ Solve the regularized linear regression problem

$$\hat{f}_h = \arg \min_{f \in \mathcal{H}_\kappa} \sum_{t=1}^k (V_{h+1}(s_{h+1,t}) - f(s_{h,t}, a_{h,t}))^2 + \lambda \|f\|_{\mathcal{H}_\kappa}^2$$

- ▶ That provides a prediction

$$\begin{aligned} \widehat{[P_h V_{h+1,k}]}(s, a) &= \widehat{f}_h(s, a) = \boldsymbol{\kappa}_{h,k}^\top(s, a)(\mathbf{K}_{h,k} + \lambda I)^{-1}\mathbf{v}_{h,k} \\ \mathbf{v}_{h,k} &= [V_{h+1}(s_{h+1,1}), V_{h+1}(s_{h+1,2}), \dots, V_{h+1}(s_{h+1,k})] \end{aligned}$$

- ▶ Also, an uncertainty quantification (variance)

$$\sigma_{h,k}^2(s, a) = \boldsymbol{\kappa}((s, a), (s, a)) - \boldsymbol{\kappa}_{h,k}^\top(s, a)(\mathbf{K}_{h,k} + \lambda I)^{-1}\boldsymbol{\kappa}_{h,k}(s, a)$$

Kernel-based optimistic value iteration (KOVI)

The prediction and variance give us an upper confidence bound on Q^* :

$$Q_{h,k}(s, a) = r_h(s, a) + \widehat{[P_h V_{h+1}]}(s, a) + \beta(\delta) \sigma_{h,k}(s, a)$$

This is then used to compute an UCB on V^* as

$$V_{h,k}(s) = \max_a Q_{h,k}(s, a)$$

Performance guarantees

Theorem [Yang et al., 2020] The regret of KOVI satisfies
 $\text{Regret}(K) = \tilde{\mathcal{O}}(H^2 \sqrt{(\mathfrak{D}^2 + \mathfrak{D} \log \mathcal{N}_\epsilon) K}).$

Performance guarantees

Theorem [Yang et al., 2020] The regret of KOVI satisfies
 $\text{Regret}(K) = \tilde{\mathcal{O}}(H^2 \sqrt{(\mathfrak{D}^2 + \mathfrak{D} \log \mathcal{N}_\epsilon) K}).$

Proof ideas:

- Prove confidence bounds

$$|\widehat{[P_h V_{h+1,k}]}(s, a) - [P_h V_{h+1,k}](s, a)| \leq \beta(\delta) \sigma_{h,k}(s, a).$$

- Using standard techniques, show

$$\text{Regret}(K) \lesssim \sum_{h,k} \beta(\delta) \sigma_{h,k}(s_{h,k}, a_{h,k}).$$

- Kernelized elliptical potential lemma (e.g., Srinivas et al., 2010)

$$\sum_{h,k} \sigma_{h,k}(s_{h,k}, a_{h,k}) \lesssim H \sqrt{K \mathfrak{D} \log(K)}.$$

Kernel-based Setting - Analysis

- We need a confidence bound of the form

$$\left| \hat{f}(s, a) - [P_h V_{h+1,k}](s, a) \right| \leq \beta(\delta) \sigma_h(s, a).$$

Kernel-based Setting - Analysis

- We need a confidence bound of the form

$$\left| \hat{f}(s, a) - [P_h V_{h+1,k}](s, a) \right| \leq \beta(\delta) \sigma_h(s, a).$$

- For a fixed $f \in \mathcal{H}_\kappa$ with non-adaptive inputs z_1, \dots, z_k ,

$$\beta(\delta) \approx \|f\|_{\mathcal{H}_\kappa} + \frac{H}{\sqrt{\lambda}} \sqrt{d \log\left(\frac{T}{\delta}\right)}$$

Kernel-based Setting - Analysis

- We need a confidence bound of the form

$$\left| \hat{f}(s, a) - [P_h V_{h+1,k}](s, a) \right| \leq \beta(\delta) \sigma_h(s, a).$$

- For a fixed $f \in \mathcal{H}_\kappa$ with non-adaptive inputs z_1, \dots, z_k ,

$$\beta(\delta) \approx \|f\|_{\mathcal{H}_\kappa} + \frac{H}{\sqrt{\lambda}} \sqrt{d \log(\frac{T}{\delta})}$$

Challenge 1: Inputs $(s_1, a_1), \dots, (s_k, a_k)$ are adaptive!

Kernel-based Setting - Analysis

- We need a confidence bound of the form

$$\left| \hat{f}(s, a) - [P_h V_{h+1, k}](s, a) \right| \leq \beta(\delta) \sigma_h(s, a).$$

- For a fixed $f \in \mathcal{H}_\kappa$ with non-adaptive inputs z_1, \dots, z_k ,

$$\beta(\delta) \approx \|f\|_{\mathcal{H}_\kappa} + \frac{H}{\sqrt{\lambda}} \sqrt{d \log(\frac{T}{\delta})}$$

Challenge 1: Inputs $(s_1, a_1), \dots, (s_k, a_k)$ are adaptive!

Solution: Self-normalized concentration inequalities for vector-valued martingales extended to kernel setting [[Abbasi-Yadkori, 2013](#), [Whitehouse et al., 2023](#)]:

$$\beta(\delta) \approx \|f\|_{\mathcal{H}_\kappa} + \frac{H}{\sqrt{\lambda}} \sqrt{\mathfrak{D} + \log(\frac{1}{\delta})}$$

Kernel-based Setting - Analysis

- We need a confidence bound of the form

$$\left| \hat{f}(s, a) - [P_h V_{h+1, k}](s, a) \right| \leq \beta(\delta) \sigma_h(s, a).$$

- For a fixed $f \in \mathcal{H}_\kappa$ with non-adaptive inputs z_1, \dots, z_k ,

$$\beta(\delta) \approx \|f\|_{\mathcal{H}_\kappa} + \frac{H}{\sqrt{\lambda}} \sqrt{d \log(\frac{T}{\delta})}$$

Challenge 1: Inputs $(s_1, a_1), \dots, (s_k, a_k)$ are adaptive!

Solution: Self-normalized concentration inequalities for vector-valued martingales extended to kernel setting [[Abbasi-Yadkori, 2013](#), [Whitehouse et al., 2023](#)]:

$$\beta(\delta) \approx \|f\|_{\mathcal{H}_\kappa} + \frac{H}{\sqrt{\lambda}} \sqrt{\mathfrak{D} + \log(\frac{1}{\delta})}$$

Challenge 2: $f = P_h V_{h+1, k}$ is **not** fixed, but depends on past data!

Kernel-based Setting - Analysis

- We need a confidence bound of the form

$$\left| \hat{f}(s, a) - [P_h V_{h+1,k}](s, a) \right| \leq \beta(\delta) \sigma_h(s, a).$$

- For a fixed $f \in \mathcal{H}_\kappa$ with non-adaptive inputs z_1, \dots, z_k ,

$$\beta(\delta) \approx \|f\|_{\mathcal{H}_\kappa} + \frac{H}{\sqrt{\lambda}} \sqrt{d \log\left(\frac{T}{\delta}\right)}$$

Challenge 1: Inputs $(s_1, a_1), \dots, (s_k, a_k)$ are adaptive!

Solution: Self-normalized concentration inequalities for vector-valued martingales extended to kernel setting [[Abbasi-Yadkori, 2013](#), [Whitehouse et al., 2023](#)]:

$$\beta(\delta) \approx \|f\|_{\mathcal{H}_\kappa} + \frac{H}{\sqrt{\lambda}} \sqrt{\mathfrak{D} + \log\left(\frac{1}{\delta}\right)}$$

Challenge 2: $f = P_h V_{h+1,k}$ is **not** fixed, but depends on past data!

Solution: Covering number argument

Covering Number Argument

- ▶ Notice that all value functions $V_{h,k}$ belong to the function class

$$\mathcal{V} = \left\{ V(s) = \min\{H, \max_a \hat{f}(s, a) + \beta\sigma(s, a)\} \right\}$$

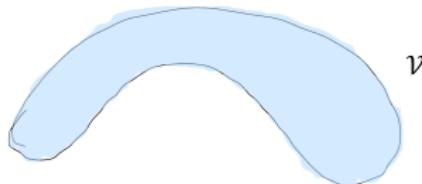
- ▶ How many functions V are required to cover \mathcal{V} up to ϵ error?

Covering Number Argument

- ▶ Notice that all value functions $V_{h,k}$ belong to the function class

$$\mathcal{V} = \left\{ V(s) = \min\{H, \max_a \hat{f}(s, a) + \beta\sigma(s, a)\} \right\}$$

- ▶ How many functions V are required to cover \mathcal{V} up to ϵ error?

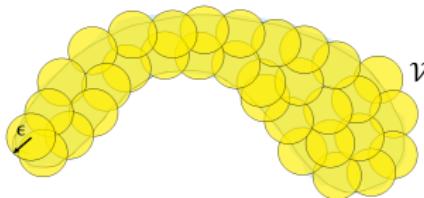


Covering Number Argument

- ▶ Notice that all value functions $V_{h,k}$ belong to the function class

$$\mathcal{V} = \left\{ V(s) = \min\{H, \max_a \hat{f}(s, a) + \beta\sigma(s, a)\} \right\}$$

- ▶ How many functions V are required to cover \mathcal{V} up to ϵ error?



Covering Number Argument

- We can now use a union-bound argument

$$\beta(\delta) = \bar{\beta}(\delta/\mathcal{N}_\epsilon) \approx \|f\|_{\mathcal{H}_k} + \frac{H}{\sqrt{\lambda}} \sqrt{\mathfrak{D} + \log \mathcal{N}_\epsilon + \frac{1}{\delta}}$$

Covering Number Argument

- We can now use a union-bound argument

$$\beta(\delta) = \bar{\beta}(\delta/\mathcal{N}_\epsilon) \approx \|f\|_{\mathcal{H}_k} + \frac{H}{\sqrt{\lambda}} \sqrt{\mathfrak{D} + \log \mathcal{N}_\epsilon + \frac{1}{\delta}}$$

- Regret ($\epsilon \approx \frac{1}{K}$) [Yang et al., 2020]

$$\text{Regret}(K) = \tilde{\mathcal{O}}(H^2 \sqrt{\mathfrak{D}^2 K + \mathfrak{D} \log \mathcal{N}_\epsilon K})$$

Covering Number Argument

- ▶ We can now use a union-bound argument

$$\beta(\delta) = \bar{\beta}(\delta/\mathcal{N}_\epsilon) \approx \|f\|_{\mathcal{H}_k} + \frac{H}{\sqrt{\lambda}} \sqrt{\mathfrak{D} + \log \mathcal{N}_\epsilon + \frac{1}{\delta}}$$

- ▶ Regret ($\epsilon \approx \frac{1}{K}$) [Yang et al., 2020]

$$\text{Regret}(K) = \tilde{\mathcal{O}}(H^2 \sqrt{\mathfrak{D}^2 K + \mathfrak{D} \log \mathcal{N}_\epsilon K})$$

- ▶ Sample Complexity

- Very smooth kernels \mathfrak{D} and $\log(\mathcal{N}_\epsilon) \approx \text{poly log}(K)$

$$\tilde{\mathcal{O}}\left(\frac{1}{\epsilon^2}\right)$$

- In general could be vacuous!

Optimistic Closure

Chowdhury and Oliveira [2023] Optimistic Closure Assumption:

$$\mathcal{V} \in \mathcal{H}_{\kappa'}$$

$$\mathcal{V} = \left\{ V(s) = \min\{H, \max_a \hat{f}(s, a) + \beta\sigma(s, a)\} \right\}$$

Optimistic Closure

Chowdhury and Oliveira [2023] Optimistic Closure Assumption:

$$\mathcal{V} \in \mathcal{H}_{\kappa'}$$

$$\mathcal{V} = \left\{ V(s) = \min\{H, \max_a \hat{f}(s, a) + \beta\sigma(s, a)\} \right\}$$

Idea: Leverage kernel mean embedding

$$\text{Regret}(K) = \tilde{\mathcal{O}}(H^2 \mathfrak{D} \sqrt{K})$$

Optimistic Closure

Chowdhury and Oliveira [2023] Optimistic Closure Assumption:

$$\mathcal{V} \in \mathcal{H}_{\kappa'}$$

$$\mathcal{V} = \left\{ V(s) = \min\{H, \max_a \hat{f}(s, a) + \beta\sigma(s, a)\} \right\}$$

Idea: Leverage kernel mean embedding

$$\text{Regret}(K) = \tilde{\mathcal{O}}(H^2 \mathfrak{D} \sqrt{K})$$

Does not hold in the linear setting!

Open Problem Vakili [2024]

- (a) Can a **no-regret** learning algorithm be designed?
- (b) What is the minimum regret growth rate with K (and also H)? And, can a learning algorithm be designed to achieve order optimal (or near-optimal) regret performance, closely aligning with the established lower bound?

Some References

- Chowdhury and Gopalan [2019]
- Yang et al. [2020]
- Domingues et al. [2021]
- Vakili and Olkhovskaya [2023]
- ...

References I

- Y. Abbasi-Yadkori. Online learning for linearly parametrized control problems. *PhD Thesis, University of Alberta*, 2013.
- Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in Neural Information Processing Systems*, 24, 2011.
- A. Ayoub, Z. Jia, C. Szepesvari, M. Wang, and L. Yang. Model-based reinforcement learning with value-targeted regression. In *International Conference on Machine Learning*, pages 463–474. PMLR, 2020.
- M. G. Azar, I. Osband, and R. Munos. Minimax regret bounds for reinforcement learning. In *International conference on machine learning*, pages 263–272. PMLR, 2017.
- S. R. Chowdhury and A. Gopalan. Online learning in kernelized Markov decision processes. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3197–3205. PMLR, 2019.
- S. R. Chowdhury and R. Oliveira. Value function approximations via kernel embeddings for no-regret reinforcement learning. In *Asian Conference on Machine Learning*, pages 249–264. PMLR, 2023.
- O. D. Domingues, P. Ménard, M. Pirotta, E. Kaufmann, and M. Valko. Kernel-based reinforcement learning: A finite-time analysis. In *International Conference on Machine Learning*, pages 2783–2792. PMLR, 2021.

References II

- S. Du, S. Kakade, J. Lee, S. Lovett, G. Mahajan, W. Sun, and R. Wang. Bilinear classes: A structural framework for provable generalization in RL. In *International Conference on Machine Learning*, pages 2826–2836. PMLR, 2021.
- C. Jin, Z. Allen-Zhu, S. Bubeck, and M. I. Jordan. Is Q-learning provably efficient? *Advances in neural information processing systems*, 31, 2018.
- C. Jin, Z. Yang, Z. Wang, and M. I. Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR, 2020.
- C. Jin, Q. Liu, and S. Miryoosefi. Bellman eluder dimension: New rich classes of RL problems, and sample-efficient algorithms. *Advances in neural information processing systems*, 34:13406–13418, 2021.
- C. Jin, Z. Yang, Z. Wang, and M. I. Jordan. Provably efficient reinforcement learning with linear function approximation. *Mathematics of Operations Research*, 48(3):1496–1521, 2023.
- T. Lattimore and C. Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- J. Lee and M.-h. Oh. Demystifying linear mdps and novel dynamics aggregation framework. In *The Twelfth International Conference on Learning Representations*, 2024.

References III

- A. Modi, J. Chen, A. Krishnamurthy, N. Jiang, and A. Agarwal. Model-free representation learning and exploration in low-rank mdps. *Journal of Machine Learning Research*, 25(6):1–76, 2024.
- A. Moulin and G. Neu. Optimistic planning by regularized dynamic programming. In *International Conference on Machine Learning*, pages 25337–25357. PMLR, 2023.
- G. Neu and C. Pike-Burke. A unifying view of optimism in episodic reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1392–1403, 2020.
- N. Srinivas, A. Krause, S. Kakade, and M. Seeger. Gaussian process optimization in the bandit setting: no regret and experimental design. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, pages 1015–1022, 2010.
- S. Vakili. Open problem: Order optimal regret bounds for kernel-based reinforcement learning. In *The Thirty Seventh Annual Conference on Learning Theory*, pages 5340–5344. PMLR, 2024.
- S. Vakili and J. Olkhovskaya. Kernelized reinforcement learning with order optimal regret bounds. *Advances in Neural Information Processing Systems*, 36, 2023.
- S. Vakili, K. Khezeli, and V. Picheny. On information gain and regret bounds in Gaussian process bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 82–90. PMLR, 2021.

References IV

- G. Weisz, P. Amortila, and C. Szepesvári. Exponential lower bounds for planning in mdps with linearly-realizable optimal action-value functions. In *Algorithmic Learning Theory*, pages 1237–1264. PMLR, 2021.
- G. Weisz, A. György, T. Kozuno, and C. Szepesvári. Confident approximate policy iteration for efficient local planning in q^π -realizable mdps. *Advances in Neural Information Processing Systems*, 35:25547–25559, 2022.
- G. Weisz, A. György, and C. Szepesvári. Online RL in linearly q^π -realizable MDPs is as easy as in linear MDPs if you learn what to ignore. *Advances in Neural Information Processing Systems*, 36, 2023.
- J. Whitehouse, A. Ramdas, and S. Z. Wu. On the sublinear regret of gp-ucb. *Advances in Neural Information Processing Systems*, 36, 2023.
- L. Yang and M. Wang. Sample-optimal parametric q-learning using linearly additive features. In *International conference on machine learning*, pages 6995–7004. PMLR, 2019.
- L. Yang and M. Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In *International Conference on Machine Learning*, pages 10746–10756. PMLR, 2020.

References V

- Z. Yang, C. Jin, Z. Wang, M. Wang, and M. Jordan. Provably efficient reinforcement learning with kernel and neural function approximations. *Advances in Neural Information Processing Systems*, 33:13903–13916, 2020.
- A. Zanette, D. Brandfonbrener, E. Brunskill, M. Pirotta, and A. Lazaric. Frequentist regret bounds for randomized least-squares value iteration. In *International Conference on Artificial Intelligence and Statistics*, pages 1954–1964. PMLR, 2020a.
- A. Zanette, A. Lazaric, M. Kochenderfer, and E. Brunskill. Learning near optimal policies with low inherent Bellman error. In H. D. III and A. Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 10978–10989. PMLR, 13–18 Jul 2020b.
- D. Zhou, J. He, and Q. Gu. Provably efficient reinforcement learning for discounted mdps with feature mapping. In *International Conference on Machine Learning*, pages 12793–12802. PMLR, 2021.