# Reinforcement learning (RL) and its challenges

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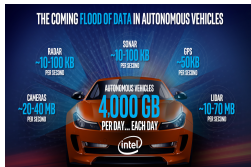


**Challenges:**

- explore or exploit: unknown or changing environments
- credit assignment problem: delayed rewards or feedback
- enormous state and action space

# Offline/Batch RL motivation: sample efficiency

- Having stored tons of history data
- Collecting new data might be expensive or time-consuming
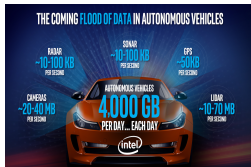


medical records



data of self-driving



clicking times of ads

- Having stored tons of history data
- Collecting new data might be expensive or time-consuming
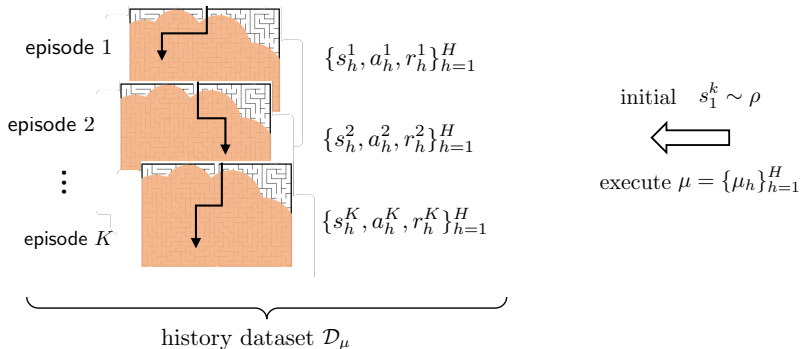


medical records



data of self-driving



clicking times of ads

**Can we design sample-efficient algorithms based on only history data?**

Given a history dataset $\mathcal{D}_\mu$ of $K$ episodes, each consisting of $H$ steps:

$$\mathcal{D}_\mu := \left\{ \left( s_1^k, a_1^k, r_1^k, \cdots, s_H^k, a_H^k, r_H^k \right) \right\}_{k=1}^{K}$$



episode 1    $\{s_h^1, a_h^1, r_h^1\}_{h=1}^{H}$

episode 2    $\{s_h^2, a_h^2, r_h^2\}_{h=1}^{H}$

$\vdots$

episode $K$    $\{s_h^K, a_h^K, r_h^K\}_{h=1}^{H}$

history dataset $\mathcal{D}_\mu$

initial   $s_1^k \sim \rho$
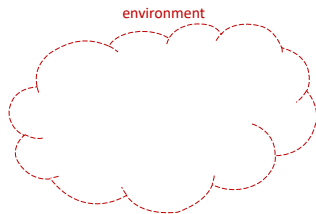
$\Longleftarrow$

execute $\mu = \{\mu_h\}_{h=1}^{H}$

**Performance metric:** Given initial state distribution $\rho$ and any accuracy level $\epsilon$. An $\epsilon$-optimal policy $\widehat{\pi} = \{\widehat{\pi}_h\}_{h=1}^H$ obeys
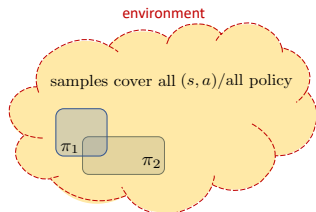
$$V_1^\star(\rho) - V_1^{\widehat{\pi}}(\rho) \le \epsilon$$

> **Goal: find an $\epsilon$-optimal policy using only history dataset**

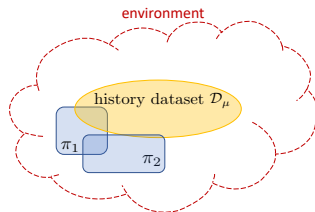— in a sample-efficient manner

environment

Online/Vanilla Offline: Uniform coverage
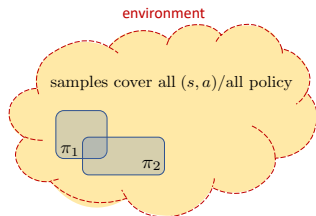
# Challenges of offline RL: partial coverage
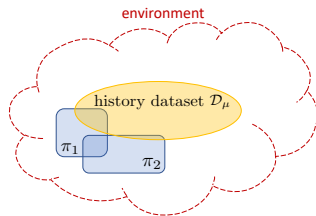


Online/Vanilla Offline: Uniform coverage

Offline RL with partial coverage

# Challenges of offline RL: partial coverage



Online/Vanilla Offline: Uniform coverage
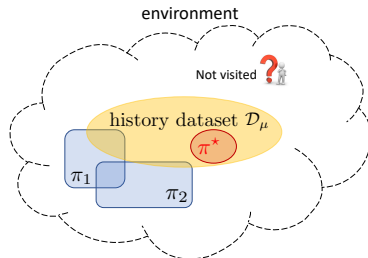
Offline RL with partial coverage

**Assumption on $\mathcal{D}_\mu$: finite single-policy concentrability**

# Key idea: pessimism/conservatism

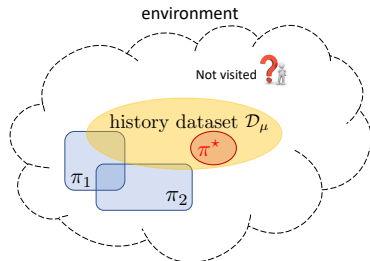*— Jin et al. '20, Rashidinejad et al. '21, Xie et al. '21*
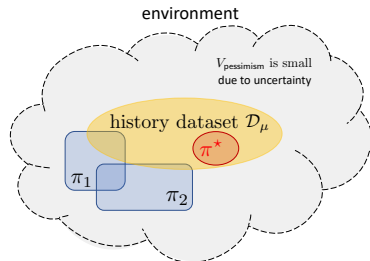


Challenge: partial coverage

# Key idea: pessimism/conservatism

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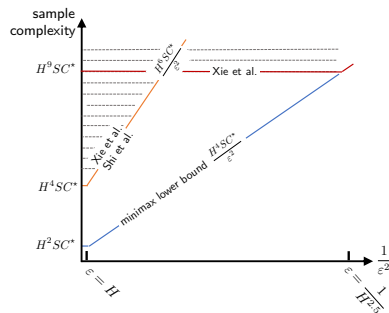


Challenge: partial coverage

Pessimism helps: $V_{\mathsf{pessimism}} < V^\star$
$V_{\mathsf{pessimism}}$ near $V^\star$

**Pessimism in Offline RL:**
add $(s, a)$-dependent penalties to reduce uncertainty damage.

# Prior art: Xie et al. '21

**Sample complexity** $T = KH = |\mathcal{D}_\mu|$

| Algorithm | Type | Sample complexity |
|---|---|---|
| VI-LCB (Xie et al., 2021) | model-based | $H^6 SC^\star/\epsilon^2$ |
| PEVI-Adv (Xie et al., 2021) | model-based | $H^4 SC^\star/\epsilon^2$ |
| lower bound (Xie et al., 2021) | n/a | $H^4 SC^\star/\epsilon^2$ |

**Sample complexity** $T = KH = |\mathcal{D}_\mu|$

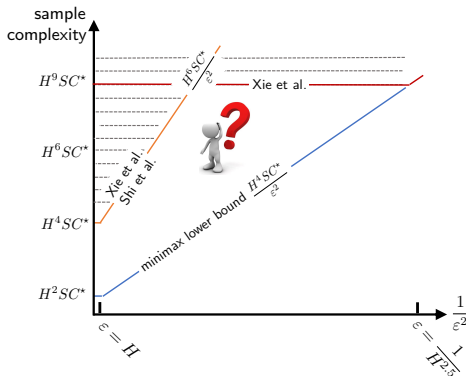| Algorithm | Type | Sample complexity |
|---|---|---|
| VI-LCB (Xie et al., 2021) | model-based | $H^6 SC^\star/\epsilon^2$ |
| PEVI-Adv (Xie et al., 2021) | model-based | $H^4 SC^\star/\epsilon^2$ |
| lower bound (Xie et al., 2021) | n/a | $H^4 SC^\star/\epsilon^2$ |



Model-based RL achieve optimal sample complexity when the accuracy level is small enough ($\epsilon \leq \frac{1}{H^{2.5}}$)
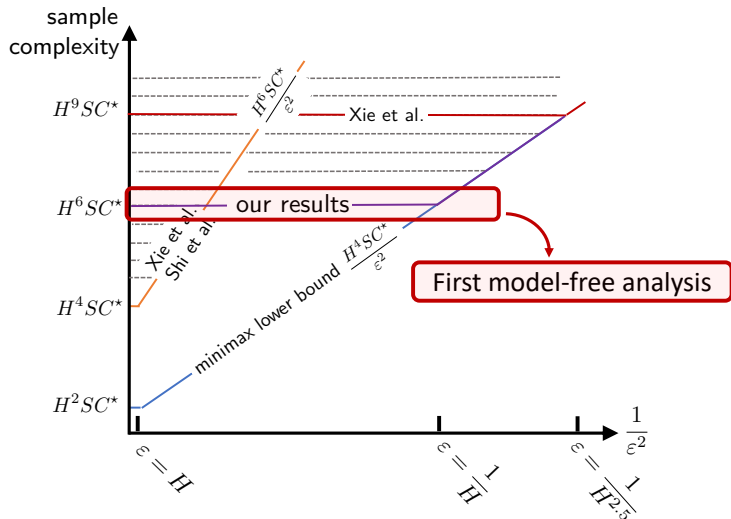
7

# No model-free offline RL analysis

# No model-free offline RL analysis



Will flexible model-free RL work?
Can we enlarge the range of accuracy level $\epsilon$?

# Our algorithm: LCB-Q-Advantage

**Theorem (Shi, Li, Wei, Chen, Chi, 2022)**

*With high prob., for $\epsilon \in \left(0, \frac{1}{H}\right]$,* LCB-Q-Advantage *can find an $\epsilon$-optimal policy $\widehat{\pi}$ as long as (up to log factor)*

$$T \gtrsim O\left(\frac{H^4 S C^\star}{\epsilon^2}\right).$$

# Our algorithm: LCB-Q-Advantage
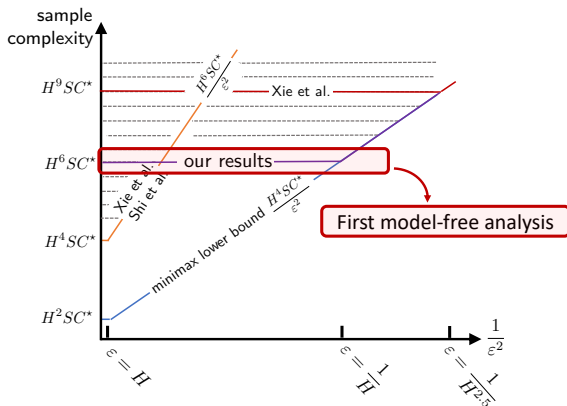
**Theorem (Shi, Li, Wei, Chen, Chi, 2022)**

*With high prob., for $\epsilon \in \left(0, \frac{1}{H}\right]$, LCB-Q-Advantage can find an $\epsilon$-optimal policy $\widehat{\pi}$ as long as (up to log factor)*

$$T \gtrsim O\left(\frac{H^4 S C^\star}{\epsilon^2}\right).$$

- model-free RL achieves optimal sample complexity for certain accuracy $\epsilon$

- optimal in a <span style="color:red">larger accuracy</span> range (improved by a factor of $H^{1.5}$)

$$\underbrace{\epsilon \leq \left(0, H^{-1}\right]}_{\text{(Our LCB-Q-Advantage)}} \qquad \text{vs.} \qquad \underbrace{\epsilon \leq \left(0, H^{-2.5}\right]}_{\text{(PEVI-Adv in [Xie et al., 2021])}}$$

# Concluding remarks



Model-free RL matches the minimax-optimal sample complexity for model-based ones!

— *in a much larger range of the accuracy level*