# Robust Offline Reinforcement Learning with Linearly Structured $f$-Divergence Regularization
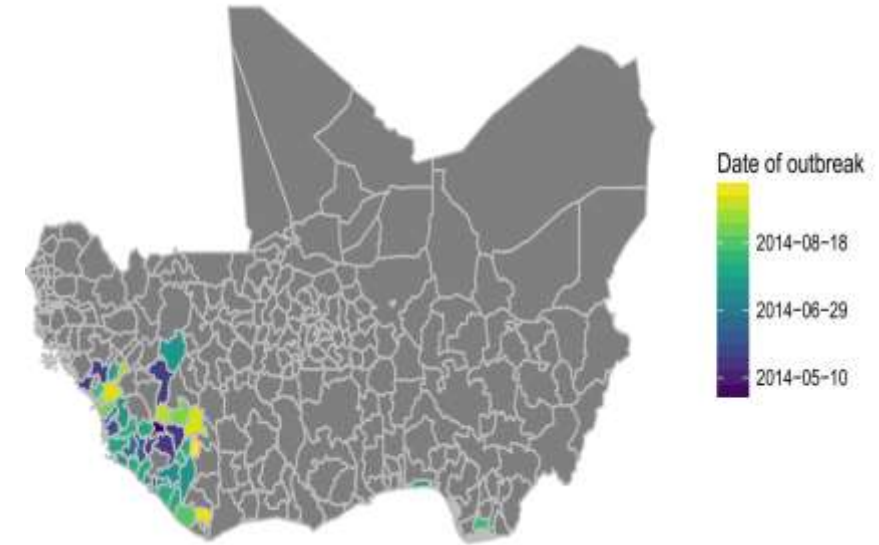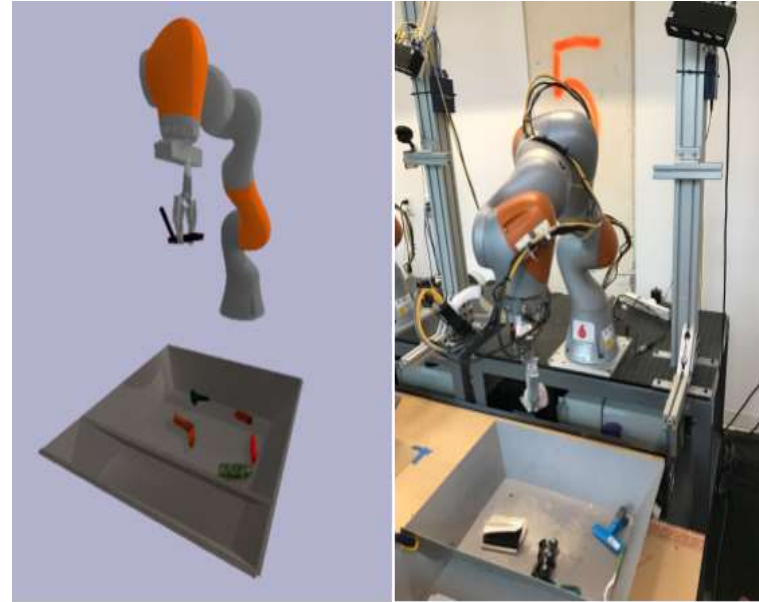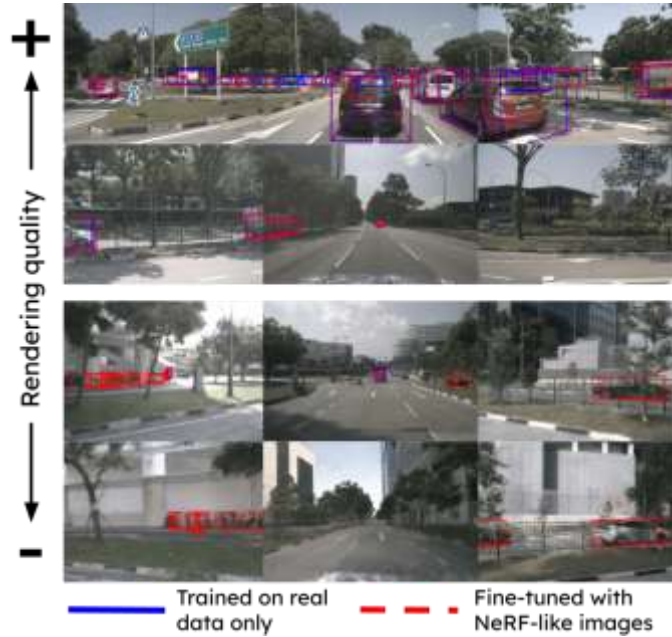
**Cheng Tang**[1], Zhishuai Liu[2], Pan Xu[2]

[1] University of Illinois Urbana-Champaign, [2] Duke University

# Content

**1** Introduction

**2** Problem Formulation

**3** Method

**4** Theoretical Analysis

**5** Experiment

# Introduction

■ Offline RL calls for **robust policy**:



● **Autonomous Driving**

● **Robotic Learning**

● **Disease Control**

[1] Lindström C, Hess G, Lilja A, et al. Are NeRFs ready for autonomous driving? Towards closing the real-to-simulation gap[C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024: 4461-4471.
[2] Bousmalis K, Levine S. Closing the simulation-to-reality gap for deep robotic learning[J]. Google Research Blog, 2017, 1.
[3] Liu Z, Clifton J, Laber E B, et al. Deep spatial q-learning for infectious disease control[J]. Journal of Agricultural, Biological and Environmental Statistics, 2023, 28(4): 749-773.

# Introduction

■ **Distributionally Robust RL**： learn more robust policy through Reinforcement Learning

- ● **d-rectangular DRMDP:** MDPs + uncertainty set

**Value function**

$$V_h^{\pi,P}(s) := \mathbb{E}^P\big[\sum_{t=h}^H r_t(s_t,a_t)\big|s_h=s,\pi\big],$$
$$Q_h^{\pi,P}(s,a) := \mathbb{E}^P\big[\sum_{t=h}^H r_t(s_t,a_t)\big|s_h=s,a_h=a,\pi\big]$$

**+**

**(s, a) - uncertainty set**

$$\mathcal{U}_{h,i}^\rho(\mu_{h,i}^0) = \{\mu : \mu \in \Delta(\mathcal{S}), D(\mu\|\mu_{h,i}^0) \le \rho\}.$$
$$\mathcal{U}_h^\rho(P_h^0) = \bigotimes_{(s,a)\in\mathcal{S}\times\mathcal{A}} \mathcal{U}_h^\rho(s,a;\boldsymbol{\mu}_h^0),$$

**‖**

**Robust value function**

$$V_h^{\pi,\rho}(s) = \inf_{P\in\mathcal{U}^\rho(P^0)} V_h^{\pi,P}(s), \quad \forall(h,s)\in[H]\times\mathcal{S}.$$

# Motivation

- **Drawbacks of d-rectangular DRMDP (<span style="color:red">hard constraint</span>):**

  - **From <span style="color:blue">theoretical</span> perspective:** need strong assumption on dual variables

  - **From <span style="color:purple">empirical</span> perspective:** solving duality problem in d-DRMDP is time-consuming

  - Existing work considers mainly TV divergence geometry, leaving blanks for cases with KL and $\chi^2$

- **RRMDP: applying regularization penalty term (<span style="color:red">soft constraint</span>) to measuring the uncertainty**

  - **From <span style="color:brown">Lagrange Duality</span> perspective:** DRMDP $\iff$ RRMDP

  - The forfeit of uncertainty set constraint makes the dual problem easier, leading to potential improvement on computation efficiency and theoretical analysis

# Content

# Problem Formulation

■ **RRMDP (Robust Regularized Markov Decision Process):** $\text{RRMDP}(S, A, H, P^0, r, \lambda, D, F)$

- Regularized robust parameter $\lambda$, probability divergence $D$, feasible set of all perturbed transition kernels F

- Regularized robust value function and Q-function:

$$V_h^{\pi,\lambda}(s) = \inf_{P \in \mathcal{F}} \mathbb{E}^P \left[ \sum_{t=h}^{H} \left[ r_t(s_t, a_t) + \lambda D(P_t(\cdot|s_t, a_t) \| P_t^0(\cdot|s_t, a_t)) \right] \Big| s_h = s, \pi \right],$$

**Penalty on divergence with nominal kernel**

$$Q_h^{\pi,\lambda}(s, a) = \inf_{P \in \mathcal{F}} \mathbb{E}^P \left[ \sum_{t=h}^{H} \left[ r_t(s_t, a_t) + \lambda D(P_t(\cdot|s_t, a_t) \| P_t^0(\cdot|s_t, a_t)) \right] \Big| s_h = s, a_h = a, \pi \right].$$

- Offline dataset and Learning goal: given $K$ trajectory $\{(s_h^\tau, a_h^\tau, r_h^\tau)\}_{h=1}^H$ and find policy $\hat{\pi}$ to minimize the robust

  Suboptimality gap: $\quad \text{SubOpt}(\hat{\pi}, s_1, \lambda) := V_1^{\star,\lambda}(s_1) - V_1^{\hat{\pi},\lambda}(s_1).$

■ **Linear MDP :**

- Known feature mapping $\phi : s \times a \to R^d$, $\sum_i \phi_i(s, a) = 1$, $\phi_i(s, a) \geq 0$

- Linear reward function and nominal transition kernel class F

$$r_h(s, a) = \langle \phi(s, a), \boldsymbol{\theta}_h \rangle, \quad P_h^0(\cdot|s, a) = \langle \phi(s, a), \boldsymbol{\mu}_h^0(\cdot) \rangle$$

# Content

# Dynamic programming principles

■ Robust regularized Bellman Equation:

$$Q_h^{\pi,\lambda}(s,a) = r_h(s,a) + \inf_{\mu_h \in \Delta(\mathcal{S})^d, P_h = \langle \phi, \mu_h \rangle} \left[ \mathbb{E}_{s' \sim P_h(\cdot|s,a)} \left[ V_{h+1}^{\pi,\lambda}(s') \right] + \lambda \langle \phi(s,a), \boldsymbol{D}(\boldsymbol{\mu}_h \| \boldsymbol{\mu}_h^0) \rangle \right],$$

$$V_h^{\pi,\lambda}(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ Q_h^{\pi,\lambda}(s,a) \right].$$

■ Existence of optimal policy

**Proposition 3.3.** Under the setting of $d$-rectangular linear RRMDP, there exists a deterministic and stationary policy $\pi^\star$, such that for any $(h,s,a) \in [H] \times \mathcal{S} \times \mathcal{A}$,

$$V_h^{\pi^\star,\lambda}(s) = V_h^{\star,\lambda}(s), \quad Q_h^{\pi^\star,\lambda}(s,a) = Q_h^{\star,\lambda}(s,a). \tag{3.7}$$

■ **Pessimism** based meta algorithm

- Step 1: estimate $w_h^\lambda$ by solving dual problem

- Step 2: construct pessimism penalty $\Gamma_h(\cdot,\cdot)$

- Step 3: compute pessimistic Q-function

**Algorithm 1** R2PVI under general $f$-divergence

**Require:** Dataset $\mathcal{D}$, Regularizer $\lambda > 0$
1: init $\hat{V}_{H+1}^\lambda(\cdot) = 0$
2: **for** episode $h = H, \cdots, 1$ **do**
3:     Compute $\Lambda_h \leftarrow \sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau)\phi(s_h^\tau, a_h^\tau)^\top + \gamma\mathbf{I}$
4:     $\hat{w}_{h,i}^\lambda(\alpha) \leftarrow \left[\Lambda_h^{-1}\left[\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau)f^*(\frac{\alpha - \hat{V}_{h+1}^\lambda(s)}{\lambda})\right]\right]^i$
        ▷ Duality Estimation for general $f$-divergence
5:     $\hat{w}_{h,i}^\lambda \leftarrow \sup_{\alpha \in \mathbb{R}}\{-\lambda\hat{w}_{h,i}^\lambda(\alpha) + \alpha\}$
6:     Construct the penalty $\Gamma_h(\cdot,\cdot)$,
7:     Estimate $\hat{Q}_h^\lambda(\cdot,\cdot) \leftarrow \min\{\langle\phi(\cdot,\cdot), \boldsymbol{\theta}_h + \hat{\boldsymbol{w}}_h^\lambda\rangle - \Gamma_h(\cdot,\cdot), H - h + 1\}^+$.
8:     Construct $\hat{\pi}_h(\cdot|\cdot) \leftarrow \arg\max_{\pi_h}\langle\hat{Q}_h^\lambda(\cdot,\cdot), \hat{\pi}_h(\cdot|\cdot)\rangle_\mathcal{A}$
    and $\hat{V}_h^\lambda(\cdot) \leftarrow \langle\hat{Q}_h^\lambda(\cdot,\cdot), \hat{\pi}_h(\cdot|\cdot)\rangle_\mathcal{A}$.
9: **end for**

# Content

# Instance-Dependent Upper Bound

■ **We provide Instance-dependent upper bound for our algorithms:**

**Theorem 5.2.** Under Assumption 3.1, for any $\delta \in (0,1)$, if we set $\gamma = 1$ and $\Gamma_h(s,a) = \beta \sum_{i=1}^{d} \|\phi_i(\cdot,\cdot)\mathbf{1}_i\|_{\Lambda_h^{-1}}$ in Algorithm 1,

- (TV) $\beta = 16Hd\sqrt{\xi_{\text{TV}}}$, where $\xi_{\text{TV}} = 2\log(1024Hd^{1/2}K^2/\delta)$;

- (KL) $\beta = 16d\lambda e^{H/\lambda}\sqrt{(H/\lambda + \xi_{\text{KL}})}$, where $\xi_{\text{KL}} = \log(1024d\lambda^2 K^3 H/\delta)$;

- ($\chi^2$) $\beta = 8dH^2(1 + 1/\lambda)\sqrt{\xi_{\chi^2}}$, where $\xi_{\chi^2} = \log(192K^5 H^6 d^3(1 + H/2\lambda)^3/\delta)$,

then with probability at least $1 - \delta$, for all $s \in \mathcal{S}$, the suboptimality of Algorithm 1 satisfies:

$$\text{SubOpt}(\hat{\pi}, s, \lambda) \leq 2\beta \left[ \sup_{P \in \mathcal{U}^\lambda(P^0)} \sum_{h=1}^{H} \mathbb{E}^{\pi^*, P} \left[ \sum_{i=1}^{d} \|\phi_i(s,a)\mathbf{1}_i\|_{\Lambda_h^{-1}} \Big| s_1 = s \right] \right].$$

$\Phi(\Lambda_h^{-1}, s)$: uncertainty function

● The upper bound relies on a novel uncertainty function

■ We further establish **information-theoretic lower bound** to illustrate the necessity of $\Phi(\Lambda_h^{-1}, s)$

■ **Comparison of the Suboptimality gap with dataset coverage**

| Algorithm | Setting | Divergence | Coverage | Suboptimality Gap |
|---|---|---|---|---|
| DRPVI (Liu & Xu, 2024b) | $d$-DRMDP | TV | full | $\tilde{O}(dH^2K^{-1/2})$ |
| DROP (Wang et al., 2024a) | $d$-DRMDP | TV | robust partial | $\tilde{O}(d^{3/2}H^2K^{-1/2})$ |
| P2MPO (TV) (Blanchet et al., 2024) | $d$-DRMDP | TV | robust partial | $\tilde{O}(d^2H^2K^{-1/2})$ |
| R2PVI-TV (**ours**) | $d$-RRMDP | TV | regularized partial | $\tilde{O}(d^2H^2K^{-1/2})$ |
| DRVI-L (Ma et al., 2022) | $d$-DRMDP | KL | robust partial | $\tilde{O}(\sqrt{\underline{\beta}}e^{H/\underline{\beta}}d^2H^{3/2}K^{-1/2})^\star$ |
| P2MPO (KL) (Blanchet et al., 2024) | $d$-DRMDP | KL | robust partial | $\tilde{O}(e^{H/\underline{\beta}}d^2H^2\rho^{-1}K^{-1/2})^\star$ |
| R2PVI-KL (**ours**) | $d$-RRMDP | KL | regularized partial | $\tilde{O}(\sqrt{\lambda}e^{H/\lambda}d^2H^{3/2}K^{-1/2})$ |
| R2PVI-$\chi^2$ (**ours**) | $d$-RRMDP | $\chi^2$ | regularized partial | $\tilde{O}(d^2H^3(1+\lambda^{-1})K^{-1/2})$ |

*\* The ⋆ denotes that the result requires an additional assumption on the KL dual variable, which is not required in **R2PVI***

- For TV divergence, R2PVI achieves <span style="color:red">nearly same</span> suboptimality gap with SOTA

- For KL divergence, R2PVI <span style="color:red">needs no extra assumption</span> to guarantee the closeness form solution

- For $\chi^2$ divergence, we are <span style="color:red">the first</span> to give theoretical guarantee under linear MDP setting

# Content

1 Introduction

2 Problem Formulation

3 Method

4 Theoretical Analysis

5 Experiment

# Experiment

■ **We want to explore:**

- The robustness of R2PVI when facing adversarial dynamics

- The role of regularizer $\lambda$ in determining the robustness of R2PVI

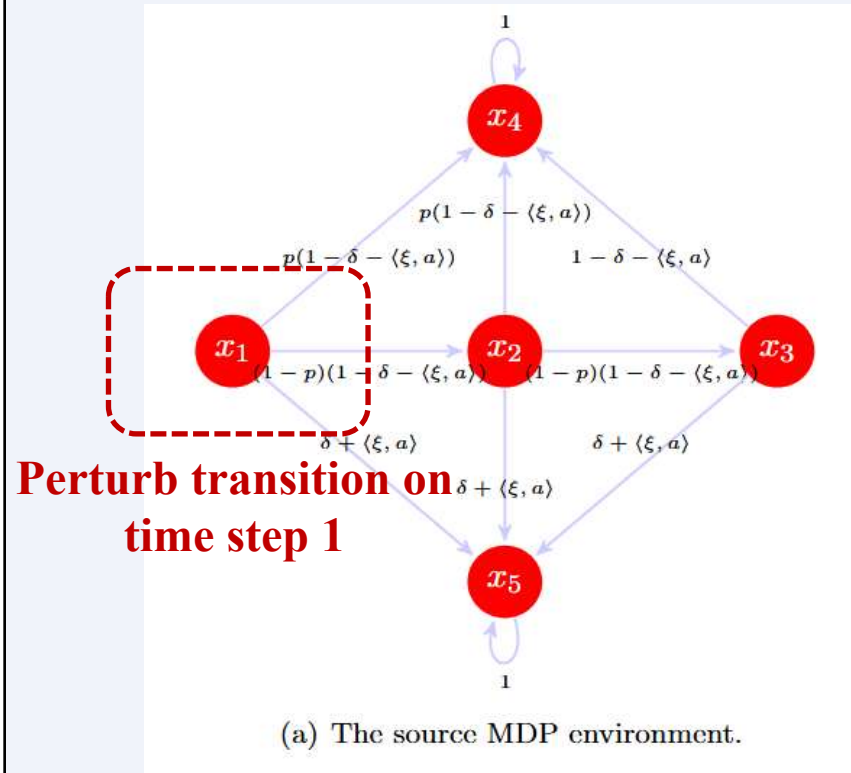- The computation cost of R2PVI compared to other robust algorithms

■ **Baselines**

| Method | PEVI | DRPVI | DRVI-L | R2PVI (ours) |
|--------|------|-------|--------|--------------|
| Framework | MDP | d-DRMDP | d-DRMDP | d-RRMDP |
| Divergence | / | TV | KL | TV/KL/$\chi^2$ |

*\* We don't compare DROP and P2MPO mentioned in the upper bound due to the lack of experiment and code base in such works.*
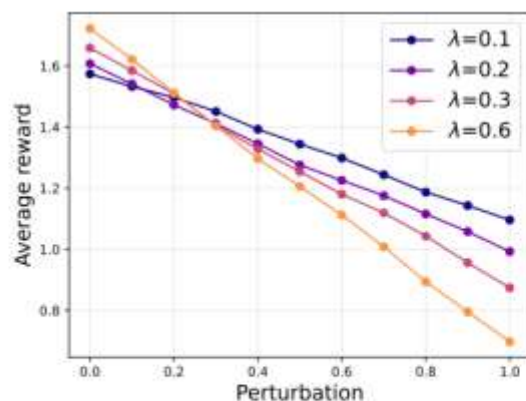
# Experiment

## ■ Task settings

### ● <u>Simulated Linear MDP</u>



(a) The source MDP environment.

**Perturb transition on time step 1**

### ● <u>American Put Option</u>



**Price fluctuates**
**(through Bernoulli Distribution)**

**Buy**    **Not Buy**

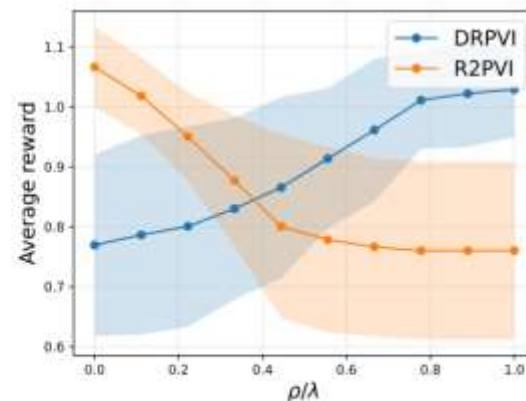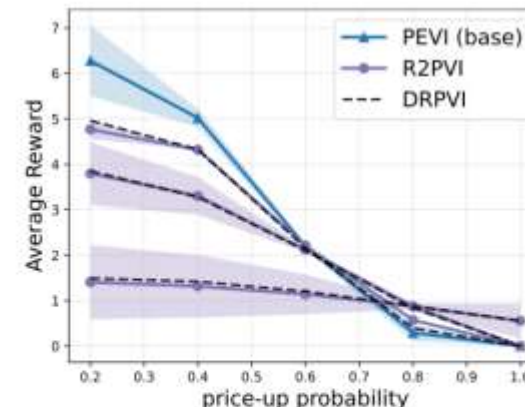**Agent**                **Dataset**
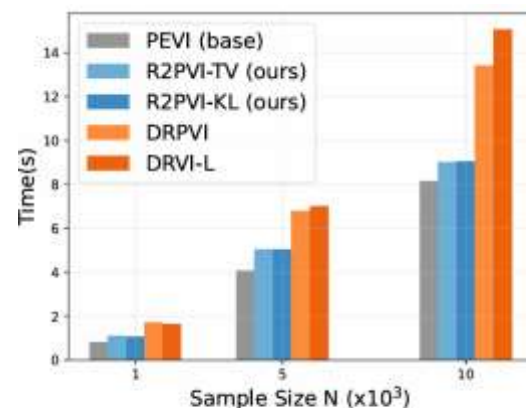
## ■ Evaluation



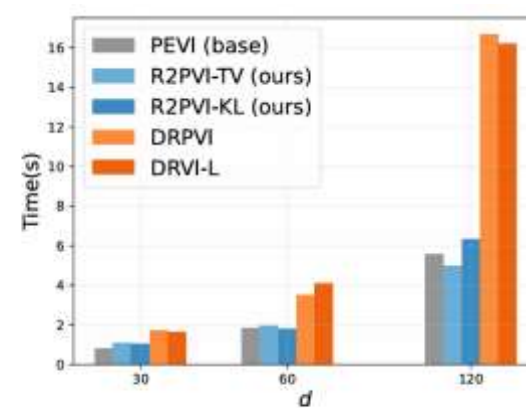(a) $\lambda = 0.1$  (b) R2PVI  (c) $q = 0.9$  (d) (up to down) $\lambda = 1, 3, 5, \rho = 0.2, 0.1, 0.025.$

(a) execution time w.r.t N.  (b) execution time w.r.t d.