# **Robust Offline Reinforcement Learning with Linearly Structured Regularization and** *f***-Divergence**

Speaker: Cheng Tang

Nov 23th, 2024

Tsinghua University



# Introduction 1 **Problem Formulation** 2 Method 3 **Theoretical Analysis** 4 Experiment 5 Conclusion 6

(2)

# Introduction



Date of outbreak

2014-08-18

2014-06-29

2014-05-10

## Sim-to-real gap:











[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.

**Cheng Tang** 

**R2PVI** 

# Introduction



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

• d-rectangular DRMDP: model the uncertainty in the dynamics and aim to achieve the

best performance under the most adversarial dynamics.



• The robust value function can be seen as the **worst** value function in an uncertainty set determined by probability divergence *D* 

# Motivation



- Drawbacks of d-rectangular DRMDP:
  - Needs strong assumption on dual variables
  - Existing algorithms rely on approximation to substitute the duality optimization, making it time consuming
  - Existing work consider mainly TV divergence geometry, leaving blanks for cases with KL and  $\chi^2$
- **RRMDP:** applying regularization penalty term measuring the uncertainty
  - Related work shows that the robust policy under RRMDP can be equivalent to DRMDP
  - The forfeit of uncertainty set constraint makes the dual problem easier, leading to potential improvement on computation efficiency and theoretical analysis

Content





(6)



## Offline MDP (Markov Decision Process): MDP(S, A, H, P<sup>0</sup>, r)

- State *s*, action *a*, reward  $r(s, a) \in [0, 1]$  (known), nominal kernel  $P^0 = \{P_h^0\}_{h=1}^H$ ,
- Value function and Q-function:

$$V_h^{\pi}(s) = \mathbb{E}^{P^0} \bigg[ \sum_{t=h}^{H} r_t(s_t, a_t) \Big| s_h = s, \pi \bigg], \qquad Q_h^{\pi}(s, a) = \mathbb{E}^{P^0} \bigg[ \sum_{t=h}^{H} r_t(s_t, a_t) \Big| s_h = s, a_h = a, \pi \bigg].$$

• 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 Suboptimality gap: SubOpt $(\pi; x) = V_1^{\pi^*}(x) - V_1^{\pi}(x)$ ,



**Cheng Tang** 

**R2PVI** 

**Tsinghua University** 



#### **RRMDP** (Robust Regularized Markov Decision Process): RRMDP(S, A, H, P<sup>0</sup>, r, λ, D, F)

- Regularized robust parameter  $\lambda$ , probability divergence *D*, feasible set of all perturbed transition kernels F
- Regularized robust value function and Q-function:

$$V_{h}^{\pi,\lambda}(s) = \inf_{P \in \mathcal{F}} \mathbb{E}^{P} \left[ \sum_{t=h}^{H} \left[ r_{t}(s_{t}, a_{t}) + \frac{\lambda D(P_{t}(|s_{t}, a_{t}) \| P_{t}^{0}(\cdot|s_{t}, a_{t})) \right]}{nominal \text{ kernel}} \right], \qquad \text{Penalty on divergence with nominal kernel}$$

$$Q_{h}^{\pi,\lambda}(s,a) = \inf_{P \in \mathcal{F}} \mathbb{E}^{P} \left[ \sum_{t=h}^{H} \left[ r_{t}(s_{t}, a_{t}) + \frac{\lambda D(P_{t}(\cdot|s_{t}, a_{t}) \| P_{t}^{0}(\cdot|s_{t}, a_{t})) \right]}{s_{h} = s, a_{h} = a, \pi} \right].$$

• 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:

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

**Optimal robust regularized value function** 

Cheng lang
------------

Linear MDP (Markov Decision Process):

- Known feature mapping  $\phi: s \times a \to R^d$ ,  $\sum_i \phi_i(s, a) = 1$ ,  $\phi_i(s, a) \ge 0$
- Linear reward function and nominal transition kernel class F

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

where  $\{\boldsymbol{\theta}_h\}_{h=1}^H$  are known vectors with bounded norm  $\|\boldsymbol{\theta}_h\|_2 \leq \sqrt{d}$  and  $\{\boldsymbol{\mu}_h^0\}_{h=1}^H$  are unknown probability measure vectors over  $\mathcal{S}$ , i.e.,  $\boldsymbol{\mu}_h^0 = (\mu_{h,1}^0, \mu_{h,2}^0, \cdots, \mu_{h,d}^0), \ \mu_{h,i}^0 \in \Delta(\mathcal{S}), \forall i \in [d].$ 

## ■ Offline d-rectangular linear RRMDP (d-RRMDP)

- Regularized robust value function and Q-function (under linear setting):  $V_{h}^{\pi,\lambda}(s) = \inf_{\mu_{t}\in\Delta(\mathcal{S})^{d}, P_{t}=\langle \phi, \mu_{t} \rangle} \mathbb{E}^{\{P_{t}\}_{t=h}^{H}} \left[ \sum_{t=h}^{H} \left[ r_{t}(s_{t}, a_{t}) + \lambda \langle \phi(s_{t}, a_{t}), \mathbf{D}(\mu_{t}||\mu_{t}^{0}) \rangle \right] \Big| s_{h} = s, \pi \right],$   $Q_{h}^{\pi,\lambda}(s,a) = \inf_{\mu_{t}\in\Delta(\mathcal{S})^{d}, P_{t}=\langle \phi, \mu_{t} \rangle} \mathbb{E}^{\{P_{t}\}_{t=h}^{H}} \left[ \sum_{t=h}^{H} \left[ r_{t}(s_{t}, a_{t}) + \lambda \langle \phi(s_{t}, a_{t}), \mathbf{D}(\mu_{t}||\mu_{t}^{0}) \rangle \right] \Big| s_{h} = s, a_{h} = a, \pi \right]$
- Optimal robust regularized value function and Q-function:

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





Robust regularized Bellman Equation:

$$Q_{h}^{\pi,\lambda}(s,a) = r_{h}(s,a) + \inf_{\substack{\mu_{h} \in \Delta(\mathcal{S})^{d}, P_{h} = \langle \phi, \mu_{h} \rangle}} \left[ \mathbb{E}_{s' \sim P_{h}(\cdot|s,a)} \left[ V_{h+1}^{\pi,\lambda}(s') \right] + \lambda \langle \phi(s,a), \mathbf{D}(\mu_{h}||\mu_{h}^{0}) \rangle \right],$$
$$V_{h}^{\pi,\lambda}(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ Q_{h}^{\pi,\lambda}(s,a) \right].$$

• This preposition shows the recurrence relationship of regularized robust value function

### Existence of optimal policy

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

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

- The existence of optimal policy guarantees the solvable of d-RRMDP with greedy policy
- The closeness of linear function class ensures the establishment of the above propositions







Experiment





# Linearity of Q-function

**Proposition 4.1.** Under Assumption 3.1, for any  $(\pi, s, a, h) \in \Pi \times S \times A \times [H]$ , we have

$$Q_{h}^{\pi,\lambda}(s,a) = \langle \phi(s,a), \theta_{h} + w_{h}^{\pi,\lambda} \rangle, \qquad (4.1)$$
  
where  $w_{h}^{\pi,\lambda} = \left(w_{h,1}^{\pi,\lambda}, w_{h,2}^{\pi,\lambda}, \cdots, w_{h,d}^{\pi,\lambda}\right)^{\top} \in \mathbb{R}^{d}$ , and  $w_{h,i}^{\pi,\lambda} = \inf_{\mu_{h,i} \in \Delta(s)} \left[\mathbb{E}^{\mu_{h,i}} \left[V_{h+1}^{\pi,\lambda}(s)\right] + \lambda D(\mu_{h,i} \| \mu_{h,i}^{0})\right].$ 

## Pessimism based algorithm

Algorithm 1 Robust Regularized Pessimistic Value Iteration (R2PVI)Step 1: estimate  $w_h^{\lambda}$  by solving dual problemAlgorithm 1 Robust Regularized Pessimistic Value Iteration (R2PVI)Require: Dataset D, Regularizer  $\lambda > 0$ 1: init  $\hat{V}_{H+1}^{\lambda}(\cdot) = 0$ 2: for episode  $h = H, \dots, 1$  do3: Compute  $\Delta_h \leftarrow \sum_{r=1}^{K} \phi(s_h^{\tau}, a_h^{\tau}) \phi(s_h^{\tau}, a_h^{\tau})^{\top} + \gamma \mathbf{I}$ 4: Obtain the parameter estimation  $\hat{w}_h^{\lambda}$ .5: Construct pessimism penalty  $\Gamma_h(\cdot, \cdot)$ 5: Construct the pessimism penalty  $\Gamma_h(\cdot, \cdot)$ 6: Estimate  $\hat{Q}_h^{\lambda}(\cdot, \cdot) \in \min(\langle \phi(\cdot, \cdot), \theta_h + w_h^{\lambda} \rangle - \Gamma_h(\cdot, \cdot), H - h + 1)^+$ .7: Construct  $\pi_h(\cdot, \cdot) \leftarrow \arg max_{\pi_h}(\hat{Q}_h^{\lambda}(\cdot, \cdot), \hat{\pi}_h(\cdot|\cdot))_{\mathcal{A}}$  and  $\hat{V}_h^{\lambda}(\cdot) \leftarrow \langle \hat{Q}_h^{\lambda}(\cdot, \cdot), \hat{\pi}_h(\cdot|\cdot)\rangle_{\mathcal{A}}$ .

#### **Cheng Tang**

#### R2PVI

#### Tsinghua University



## **Dual form of TV**

$$\inf_{\mu \in \Delta(S)} \mathbb{E}_{s \sim \mu} V(s) + \lambda D_{\mathrm{TV}}(\mu \| \mu^0) = \mathbb{E}_{s \sim \mu^0} [V(s)]_{\min_{s'}(V(s')) + \lambda}.$$

**Close form solution** 

## **Specific applicable algorithm**

• Obtain by least square regression:

$$\hat{\boldsymbol{w}}_{h}^{\lambda} = \operatorname*{argmin}_{\boldsymbol{w} \in R^{d}} \sum_{\tau=1}^{K} \left( [\hat{V}_{h+1}^{\lambda}(s_{h+1}^{\tau})]_{\alpha_{h+1}} - \boldsymbol{\phi}(s_{h}^{\tau}, a_{h}^{\tau})^{\top} \boldsymbol{w} \right)^{2} + \gamma \|\boldsymbol{w}\|_{2}^{2},$$

Algorithm 2 Robust Regularized Pessimistic Value Iteration under TV distance (R2PVI-TV  
Require: Dataset 
$$\mathcal{D}$$
, regularizer  $\lambda > 0$ ,  $\gamma > 0$  and parameter  $\beta$   
1: init  $\hat{V}_{H+1}^{\lambda}(\cdot) = 0$   
2: for episode  $h = H, \dots, 1$  do  
3:  $\Lambda_h \leftarrow \sum_{\tau=1}^{K} \phi(s_h^{\tau}, a_h^{\tau})(\phi(s_h^{\tau}, a_h^{\tau}))^{\top} + \gamma \mathbf{I}$   
4:  $\alpha_{h+1} \leftarrow \min_{s \in \mathcal{S}}(\hat{V}_{h+1}^{\lambda}(s)) + \lambda$   
5:  $\hat{w}_h^{\lambda} \leftarrow \Lambda_h^{-1}(\sum_{\tau=1}^{K} \phi(s_h^{\tau}, a_h^{\tau})[\hat{V}_{h+1}^{\lambda}(s_{h+1}^{\tau})]_{\alpha_{h+1}})$   
6:  $\overline{\Gamma}_h(\cdot, \cdot) \leftarrow \beta \sum_{i=1}^{a} \|\phi_i(\cdot, \cdot)\mathbf{1}_i\|_{\Lambda_h^{-1}}$   
7:  $\hat{Q}_h^{\lambda}(\cdot, \cdot) \leftarrow \min(\phi(\cdot, \cdot)^{\top}(\theta_h + \hat{w}_h^{\lambda}) - \Gamma_h(\cdot, \cdot), H - h + 1)^+$   
8:  $\hat{\pi}_h(\cdot|\cdot) \leftarrow \operatorname{argmax}_{\pi_h}\langle \hat{Q}_h^{\lambda}(\cdot, \cdot), \pi_h(\cdot|\cdot) \rangle_{\mathcal{A}}$  and  $\hat{V}_h^{\lambda}(\cdot) \leftarrow \langle \hat{Q}_h^{\lambda}(\cdot, \cdot), \hat{\pi}_h(\cdot|\cdot) \rangle_{\mathcal{A}}$ 

• Specifically designed penalty with  $\beta_{TV}_{9:}^{8:}$ 



#### **Dual form of KL**

$$\inf_{\mu \in \Delta(S)} \mathbb{E}_{s \sim \mu} V(s) + \lambda D_{\mathrm{KL}}(\mu \| \mu^0) = -\lambda \log \mathbb{E}_{s \sim \mu^0} \left[ e^{-\frac{V(s)}{\lambda}} \right].$$

Logarithm may generate error

## Specific applicable algorithm

- **Obtain by least square regression:** Algorithm 3 Robust Regularized Pessimistic Value Iteration under KL distance (R2PVI-KL)  $\hat{w}_h' = \operatorname*{argmin}_{w \in \mathbb{R}^d} \sum_{ au=1}^K \left( e^{-rac{\hat{v}_{h+1}^\lambda(s_{h+1}^ au)}{\lambda}} - \phi(s_h^ au, a_h^ au)^ op w 
  ight)^2 + \gamma \|w\|_2^2.$ **Require:** Dataset  $\mathcal{D}$ , regularizer  $\lambda > 0, \gamma > 0$  and parameter  $\beta$ 1: init  $\hat{V}_{H+1}^{\lambda}(\cdot) = 0$ 2: for episode  $h = H, \dots, 1$  do 3: 
  $$\begin{split} & \mathbf{\Lambda}_{h} \leftarrow \sum_{\tau=1}^{K} \phi(s_{h}^{\tau}, a_{h}^{\tau}) (\phi(s_{b}^{\tau}, a_{h}^{\tau}))^{\top} + \gamma \mathbf{I} \\ & \mathbf{\Psi}_{h}^{\prime} \leftarrow \mathbf{\Lambda}_{h}^{-1} \left( \sum_{\tau=1}^{K} \phi(s_{h}^{\tau}, a_{h}^{\tau}) e^{-\frac{\tilde{V}_{h+1}^{\lambda}(s_{h+1}^{\tau})}{\lambda}} \right) \\ & 5: \quad \hat{\mathbf{W}}_{h}^{\lambda} \leftarrow -\lambda \log \max\{ \hat{\mathbf{W}}_{h}^{\prime}, e^{-\tilde{H}/\lambda} \} \end{split}$$
  Clip the  $\widehat{W}'_h$  with lower bound o  $\mathbb{E}_{s \sim \mu^0} e^{-\widehat{V}_{h+1}^{\lambda}(s)/\lambda}$ 6:  $\Gamma_h(\cdot, \cdot) \leftarrow \beta \sum_{i=1}^d \|\phi_i(\cdot, \cdot) \mathbf{1}_i\|_{\Lambda_h^{-1}}$ 7:  $\hat{Q}_h^{\lambda}(\cdot, \cdot) \leftarrow \min(\phi(\cdot, \cdot)^{\top}(\boldsymbol{\theta}_h + \hat{w}_h') - \Gamma_h(\cdot, \cdot), H - h + 1)^+$ 8:  $\hat{\pi}_h(\cdot|\cdot) \leftarrow \operatorname{argmax}_{\pi_h} \langle \hat{Q}_h^{\lambda}(\cdot, \cdot), \pi_h(\cdot|\cdot) \rangle_{\mathcal{A}} \text{ and } \hat{V}_h^{\lambda}(\cdot) \leftarrow \langle \hat{Q}_h^{\lambda}(\cdot, \cdot), \hat{\pi}_h(\cdot|\cdot) \rangle_{\mathcal{A}}$ 9: end for
- Specifically designed penalty with  $\beta_{KL}$



# **Dual form of \chi^2**

$$\inf_{\mu \in \Delta(S)} \mathbb{E}_{s \sim \mu} V(s) + \lambda D_{\chi^2}(\mu \| \mu^0) = \left\{ \sup_{\alpha \in [V_{\min}, V_{\max}]} \left\{ \mathbb{E}_{s \sim \mu^0} [V(s)]_{\alpha} - \frac{1}{4\lambda} \operatorname{Var}_{s \sim \mu^0} [V(s)]_{\alpha} \right\} \right\}$$

Specific applicable algorithm

• estimate 
$$w_h^{\lambda}$$
 by solving dual problem  

$$\hat{\mathbb{E}}^{\mu_{h,i}^0}[\hat{V}_{h+1}^{\lambda}(s)]_{\alpha} = \left[ \operatorname{argmin}_{w \in \mathbb{R}^d} \sum_{\tau=1}^{K} ([\hat{V}_{h+1}^{\lambda}(s_{h+1}^{\tau})]_{\alpha} - \phi(s_h^{\tau}, a_h^{\tau})^{\top}w)^2 + \gamma \|w\|_2^2 \right]_{[0,H]}^i,$$

$$\hat{\mathbb{E}}^{\mu_{h,i}^0}[\hat{V}_{h+1}^{\lambda}(s)]_{\alpha}^2 = \left[ \operatorname{argmin}_{w \in \mathbb{R}^d} \sum_{\tau=1}^{K} ([\hat{V}_{h+1}^{\lambda}(s_{h+1}^{\tau})]_{\alpha}^2 - \phi(s_h^{\tau}, a_h^{\tau})^{\top}w)^2 + \gamma \|w\|_2^2 \right]_{[0,H]}^i,$$

$$\hat{v}_{h,i}^{\lambda} = \sup_{\alpha \in [(\hat{V}_{h+1}^{\lambda})\min, (\hat{V}_{h+1}^{\lambda})\max]} \left\{ \hat{\mathbb{E}}^{\mu_{h,i}^0}[\hat{V}_{h+1}^{\lambda}(s)]_{\alpha} - \frac{1}{4\lambda}\widehat{\operatorname{Var}}^{\mu_{h,i}^0}[\hat{V}_{h+1}^{\lambda}(s)]_{\alpha} \right\}$$

$$= \max_{\alpha \in [(\hat{V}_{h+1}^{\lambda})\min, (\hat{V}_{h+1}^{\lambda})\max]} \left\{ \hat{\mathbb{E}}^{\mu_{h,i}^0}[\hat{V}_{h+1}^{\lambda}(s)]_{\alpha} + \frac{1}{4\lambda} (\hat{\mathbb{E}}^{\mu_{h,i}^0}[\hat{V}_{h+1}^{\lambda}(s)]_{\alpha})^2 - \frac{1}{4\lambda}\hat{\mathbb{E}}^{\mu_{h,i}^0}[\hat{V}_{h+1}^{\lambda}(s)]_{\alpha}^2 \right\}$$

• Specifically designed penalty with  $\beta_{\chi^2}$ 

Algorithm 4 Robust Regularized Pessimistic Value Iteration under 
$$\chi^2$$
 distance (R2PVI- $\chi^2$ )Require: Dataset  $\mathcal{D}$ , regularizer  $\lambda > 0, \gamma > 0$  and parameter  $\beta$ 1: init  $\hat{V}_{H+1}^{\lambda}(\cdot) = 0$ 2: for episode  $h = H, \cdots, 1$  do3:  $\Lambda_h \leftarrow \sum_{\tau=1}^{K} \phi(s_h^{\tau}, a_h^{\tau})(\phi(s_h^{\tau}, a_h^{\tau}))^{\top} + \gamma \mathbf{I}$ 4:  $\hat{\mathbb{E}}_{h,i}^{\mu_{h,i}}[\hat{V}_{h+1}^{\lambda}(s)]_{\alpha} \leftarrow [\Lambda_h^{-1}(\sum_{\tau=1}^{K} \phi(s_h^{\tau}, a_h^{\tau})^{\top}[\hat{V}_{h+1}^{\lambda}(s_{h+1}^{\tau})]_{\alpha})]_{[0,H]}$ 4:  $\hat{\mathbb{E}}_{h,i}^{\mu_{h,i}^{0}}[\hat{V}_{h+1}^{\lambda}(s)]_{\alpha}^{2} \leftarrow [\Lambda_{h}^{-1}(\sum_{\tau=1}^{K} \phi(s_{h}^{\tau}, a_{h}^{\tau})^{\top}[\hat{V}_{h+1}^{\lambda}(s_{h+1}^{\tau})]_{\alpha}^{2})]_{[0,H^2]}$ 5:  $\hat{\mathbb{E}}_{h,i}^{\mu_{h,i}^{0}}[\hat{V}_{h+1}^{\lambda}(s)]_{\alpha}^{2} \leftarrow [\Lambda_{h}^{-1}(\sum_{\tau=1}^{K} \phi(s_{h}^{\tau}, a_{h}^{\tau})^{\top}[\hat{V}_{h+1}^{\lambda}(s_{h+1}^{\tau})]_{\alpha}^{2})]_{[0,H^2]}$ 6: Estimate  $\hat{w}_{h,i}^{\lambda}$  according to (4.9)7:  $\Gamma_h(\cdot, \cdot) \leftarrow \beta \sum_{i=1}^{d} \|\phi_i(\cdot, \cdot)\|_{\Lambda_h^{-1}}^{-1}$ 8:  $Q_h^{\lambda}(\cdot, \cdot) \leftarrow \min(\phi(\cdot, \cdot)^{\top}(\theta_h + \hat{w}_h^{\lambda}) - \Gamma_h(\cdot, \cdot), H - h + 1)^+$ 9:  $\hat{\pi}_h(\cdot|\cdot) \leftarrow \operatorname{argmax}_{\pi_h} \langle \hat{Q}_h^{\lambda}(\cdot, \cdot), \pi_h(\cdot|\cdot) \rangle_{\mathcal{A}}$  and  $\hat{V}_h^{\lambda}(\cdot) \leftarrow \langle \hat{Q}_h^{\lambda}(\cdot, \cdot), \hat{\pi}_h(\cdot|\cdot) \rangle_{\mathcal{A}}$ 10: end for





#### 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  $\Gamma_h(s,a) = \beta \sum_{i=1}^d \|\phi_i(\cdot,\cdot)\mathbf{1}_i\|_{\mathbf{\Lambda}_h^{-1}}$  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  $\xi_{\chi^2} = \log(192K^5H^6d^3(1+H/2\lambda)^3/\delta)$ ,

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

SubOpt
$$(\hat{\pi}, s, \lambda) \leq 2\beta \left[ \sup_{P \in \mathcal{U}^{\lambda}(P^{0})} \sum_{h=1}^{H} \mathbb{E}^{\pi^{*}, P} \left[ \sum_{i=1}^{d} \|\phi_{i}(s, a) \mathbf{1}_{i}\|_{\Lambda_{h}^{-1}} |s_{1} = s \right] \right].$$
  
 $\Phi(\Lambda_{h}^{-1}, s)$ : uncertainty function

• The upper bound relies on a novel uncertainty function



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

**Theorem 6.1.** Given a regularizer  $\lambda$ , dimension d, horizon length H and sample size  $K > \max\{\tilde{O}(d^6), \tilde{O}(d^3H^2/\lambda^2)\}$ , there exists a class of d-rectangular linear RRMDPs  $\mathcal{M}$  and an offline dataset  $\mathcal{D}$  of size K such that for all  $s \in S$  and any divergence D among  $D_{\text{TV}}, D_{\text{KL}}$  and  $D_{\chi^2}$ , with probability at least  $1 - \delta$ , we have  $\inf_{\hat{\pi}} \sup_{M \in \mathcal{M}} \operatorname{SubOpt}(M, \hat{\pi}, s, \lambda, D) \geq c \cdot \Phi(\Lambda_h^{-1}, s)$ , where c is a universal constant.

#### The construction of hard instance

- The nominal environment is constructed by inserting an error into the environment with two absorbing states
- The perturbed environment resembles the nominal kernel besides a controllable perturbation



(a) The nominal environment.

(b) The perturbed environment under time step h.



## Comparison of the Suboptimality gap with dataset coverage

Algorithm	Setting	Robust	Divergence	Coverage	Suboptimality gap
DRPVI	DRMDP	ρ	$\mathbf{TV}$	full	$ ilde{O}(d^{3/2}H^2K^{-1/2})$
DROP	DRMDP	ρ	$\mathbf{TV}$	robust partial	$ ilde{O}(d^{3/2}H^2K^{-1/2})$
P2MPO (TV)	DRMDP	ρ	$\mathbf{TV}$	robust partial	$ ilde{O}(d^2H^2K^{-1/2})$
R2PVI-TV	RRMDP	λ	$\mathbf{TV}$	regularized partial	$ ilde{O}(d^2H^2K^{-1/2})$
DRVI-L	DRMDP	ρ	KL	robust partial	$ ilde{O}(\sqrt{eta}e^{H/eta}d^2H^{3/2}K^{-1/2})^\star$
P2MPO (KL)	DRMDP	ρ	KL	robust partial	$ ilde{O}(e^{H/eta} d^2 H^2  ho^{-1} K^{-1/2})^{\star}$
R2PVI-KL	RRMDP	λ	KL	regularized partial	$ ilde{O}(\sqrt{\lambda}e^{H/\lambda}d^2H^{3/2}K^{-1/2})$
<b>R2PVI-</b> $\chi^2$	RRMDP	λ	$\chi^2$	regularized partial	$ ilde{O}(d^2 H^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, our algorithm achieves nearly same suboptimality gap with SOTA
- For KL divergence, our algorithm needs no extra assumption to guarantee the closeness form solution
- For  $\chi^2$  divergence, we are the first to give algorithms under linear MDP setting with  $\chi^2$  divergence







# 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.

**Cheng Tang** 

**R2PVI** 

# Experiment



# ■ Task settings





#### **Cheng Tang**

# **Simulated Linear MDP**



# Evaluation



- Compared to non robust algorithm (PEVI), R2PVI learns robust policy under all divergence measure
- Robust parameter  $\lambda$  serves as regularization to adjust the robustness of the policy
- $\lambda$  plays a similar role in d-RRMDP as inverse robust parameter  $1/\rho$  in d-DRMDP



# Evaluation



- When *N* and *d* are large, the computation cost of DRPVI and DRVI-L increase rapidly
- The computation cost of R2PVI is as low as PEVI.
- **R2PVI** can achieve the same robust performance with DRPVI and DRVI-L







# Contribution

- We propose a novel d-RRMDP framework and establish dynamic planning principles
- We derive dual formulations of Q-functions under TV, KL,  $\chi^2$  divergence, and admit their linear representations
- We design meta-algorithms, R2PVI, in our setting and provide specific applicable algorithms under TV, KL,  $\chi^2$  divergence
- We provide instance-dependent upper bounds of our algorithms with a general form  $\beta \sup_{P \in \mathcal{U}^{\lambda}(P^{0})} \sum_{h=1}^{H} \mathbb{E}^{\pi^{*},P} \Big[ \sum_{i=1}^{d} \|\phi_{i}(s,a)\mathbf{1}_{i}\|_{\Lambda_{h}^{-1}} \|s_{1} = s \Big],$

and then construct theoretical-lower bound to illustrate that the general form is intrinsic

• We conduct extensive experiment to illustrate robustness and time efficiency of our algorithms



# Thank you!

