

Distributionally Robust RL with Interactive Data Collection

Fundamental Hardness and Near-Optimal Algorithm

Miao Lu

Stanford University



INFORMS 2024

Joint work with



Han Zhong
PKU



Jose Blanchet
Stanford



Tong Zhang
UIUC

Lu, M., Zhong, H., Zhang, T., & Blanchet, J. (2024). Distributionally Robust Reinforcement Learning with Interactive Data Collection: Fundamental Hardness and Near-Optimal Algorithm. *Preliminary version at 37th Conference on Neural Information Processing Systems (NeurIPS), 2024*

Work supported by grants NSF 2118199, 2229012, 2312204 and FA9550-20-1-0397.

Outline

Background and Problem Setup

Hardness Result under Interactive Data Collection

Vanishing Minimal Value Assumption and Algorithm Design

Future Works

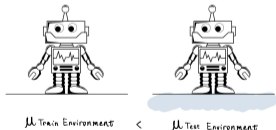
RL and Sim-to-Real Gaps

► RL achieves tremendous success in:

- Robotics
- Healthcare
- Recommendation systems
- Autonomous driving
- Training large language models

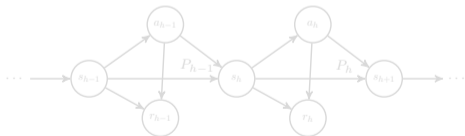
► Challenge: **sim-to-real gap**

- Training env. \neq Testing env.
- Cause degeneration of performance



Basics of RL

- MDP: $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, \{P_h\}_{h=1}^H, \{R_h\}_{h=1}^H)$.
- \mathcal{S} : state space, \mathcal{A} : action space.
- $R : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$: reward function
- H : Horizon length.
- P_h^* : transition distribution of training Env.
- P_h' : transition distribution of testing Env.



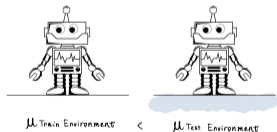
RL and Sim-to-Real Gaps

► RL achieves tremendous success in:

- Robotics
- Healthcare
- Recommendation systems
- Autonomous driving
- Training large language models

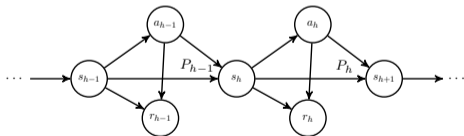
► Challenge: **sim-to-real gap**

- Training env. \neq Testing env.
- Cause degeneration of performance



Basics of RL

- MDP: $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, \{P_h\}_{h=1}^H, \{R_h\}_{h=1}^H)$.
- \mathcal{S} : state space, \mathcal{A} : action space.
- $R : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$: reward function
- H : Horizon length.
- P_h^* : transition distribution of training Env.
- P_h' : transition distribution of testing Env.



Robust Markov Decision Processes

- Goal: using data collected from training Env. $P^* = \{P_h^*\}_{h=1}^H$, find a policy π^* that

$$\pi^* := \arg \sup_{\pi \in \Pi} V_{1, P^*, \Phi}^\pi(s) := \arg \sup_{\pi \in \Pi} \inf_{\substack{P_h \in \Phi(P_h^*) \\ 1 \leq h \leq H}} \mathbb{E}_P^\pi \left[\sum_{h=1}^H R_h(s_h, a_h) \middle| s_1 = s \right].$$

- Robust set Φ : set of testing environment distributions
- This work: we focus on Total Variation (TV) distance robust set:

$$\Phi(P) = \bigotimes_{(s,a) \in \mathcal{S} \times \mathcal{A}} \mathcal{P}_\rho(s, a; P),$$

$$\mathcal{P}_\rho(s, a; P) := \left\{ \tilde{P}(\cdot) \in \Delta(\mathcal{S}) : D_{\text{TV}}(\tilde{P}(\cdot) \| P(\cdot | s, a)) \leq \rho \right\}.$$

- $\rho \in [0, 1)$ is the level of robustness.

- Robust Bellman optimal equation:

$$V_{h, P^*, \Phi}^* = \max_{a \in \mathcal{A}} Q_{h, P^*, \Phi}^*(s, a), \quad Q_{h, P^*, \Phi}^*(s, a) = R_h(s, a) + \inf_{P_h \in \Phi(P_h^*)} \mathbb{E}_P^\pi [V_{h+1, P^*, \Phi}^*].$$

Robust Markov Decision Processes

- Goal: using data collected from training Env. $P^* = \{P_h^*\}_{h=1}^H$, find a policy π^* that

$$\pi^* := \arg \sup_{\pi \in \Pi} V_{1, P^*, \Phi}^\pi(s) := \arg \sup_{\pi \in \Pi} \inf_{\substack{P_h \in \Phi(P_h^*) \\ 1 \leq h \leq H}} \mathbb{E}_P^\pi \left[\sum_{h=1}^H R_h(s_h, a_h) \middle| s_1 = s \right].$$

- Robust set Φ : set of testing environment distributions

- This work: we focus on Total Variation (TV) distance robust set:

$$\Phi(P) = \bigotimes_{(s,a) \in \mathcal{S} \times \mathcal{A}} \mathcal{P}_\rho(s, a; P),$$

$$\mathcal{P}_\rho(s, a; P) := \left\{ \tilde{P}(\cdot) \in \Delta(\mathcal{S}) : D_{\text{TV}}(\tilde{P}(\cdot) \| P(\cdot | s, a)) \leq \rho \right\}.$$

- $\rho \in [0, 1)$ is the level of robustness.

- Robust Bellman optimal equation:

$$V_{h, P^*, \Phi}^* = \max_{a \in \mathcal{A}} Q_{h, P^*, \Phi}^*(s, a), \quad Q_{h, P^*, \Phi}^*(s, a) = R_h(s, a) + \inf_{P_h \in \Phi(P_h^*)} \mathbb{E}_P^\pi [V_{h+1, P^*, \Phi}^*].$$

Robust Markov Decision Processes

- Goal: using data collected from training Env. $P^* = \{P_h^*\}_{h=1}^H$, find a policy π^* that

$$\pi^* := \arg \sup_{\pi \in \Pi} V_{1, P^*, \Phi}^\pi(s) := \arg \sup_{\pi \in \Pi} \inf_{\substack{P_h \in \Phi(P_h^*) \\ 1 \leq h \leq H}} \mathbb{E}_P^\pi \left[\sum_{h=1}^H R_h(s_h, a_h) \middle| s_1 = s \right].$$

- Robust set Φ : set of testing environment distributions
- This work: we focus on Total Variation (TV) distance robust set:

$$\Phi(P) = \bigotimes_{(s,a) \in \mathcal{S} \times \mathcal{A}} \mathcal{P}_\rho(s, a; P),$$

$$\mathcal{P}_\rho(s, a; P) := \left\{ \tilde{P}(\cdot) \in \Delta(\mathcal{S}) : D_{\text{TV}}(\tilde{P}(\cdot) \| P(\cdot | s, a)) \leq \rho \right\}.$$

- $\rho \in [0, 1)$ is the level of robustness.

- Robust Bellman optimal equation:

$$V_{h, P^*, \Phi}^* = \max_{a \in \mathcal{A}} Q_{h, P^*, \Phi}^*(s, a), \quad Q_{h, P^*, \Phi}^*(s, a) = R_h(s, a) + \inf_{P_h \in \Phi(P_h^*)} \mathbb{E}_P^\pi [V_{h+1, P^*, \Phi}^*].$$

Recap: Learn the optimal robust policy π^* using data of training Env. P^*

What kind of data do we have?

	Mechanism	Explanation
Prior Work	Generative Model	can query any (s, a, h) to obtain $s' \sim P_h^*(\cdot s, a)$
	Offline Learning	pre-collected data $\{(s_h^i, a_h^i, s_{h+1}^i)\}_{i=1}^N \sim \mu_h(s, a) \otimes P_h^*(s' s, a)$

This work	Interactive data collection	<i>interact with P^* by algorithm-dependent policy!</i>
-----------	-----------------------------	--

	Mechanism	Explanation
Prior Work	Generative Model	can query any (s, a, h) to obtain $s' \sim P_h^*(\cdot s, a)$
	Offline Learning	pre-collected data $\{(s_h^i, a_h^i, s_{h+1}^i)\}_{i=1}^N \sim \mu_h(s, a) \otimes P_h^*(s' s, a)$
This work	Interactive data collection	<i>interact with P^* by algorithm-dependent policy!</i>

Data Generation Mechanism

Interaction protocol:

- ▶ Interact with training Env. P^* for $K \in [N]$ episodes.
- ▶ In each episode $k \in [K]$, use policy π^k to collect data $\{(s_h^k, a_h^k, r_h^k, s_{h+1}^k)\}_{h=1}^H$.
- ▶ After episode k , use historical data to update policy to π^{k+1}

Metric:

- ▶ Online regret:

$$\text{Regret}_{\Phi}(K) := \sum_{k=1}^K V_{1,P^*,\Phi}^*(s_1) - V_{1,\pi^k,\Phi}(s_1).$$

- ▶ Sample complexity: # episodes of interaction suffice to learn ε -optimal robust policy, i.e.

$$V_{1,P^*,\Phi}^*(s_1) - V_{1,\hat{\pi},\Phi}(s_1) \leq \varepsilon \quad (1)$$

Question:

"Can we design a provably sample-efficient robust RL algorithm using interactive data collection in the training environment?"

Data Generation Mechanism

Interaction protocol:

- ▶ Interact with training Env. P^* for $K \in [N]$ episodes.
- ▶ In each episode $k \in [K]$, use policy π^k to collect data $\{(s_h^k, a_h^k, r_h^k, s_{h+1}^k)\}_{h=1}^H$.
- ▶ After episode k , use historical data to update policy to π^{k+1}

Metric:

- ▶ Online regret:

$$\text{Regret}_{\Phi}(K) := \sum_{k=1}^K V_{1, P^*, \Phi}^*(s_1) - V_{1, P^*, \Phi}^{\pi^k}(s_1).$$

- ▶ Sample complexity: # episodes of interaction suffice to learn ε -optimal robust policy, i.e.

$$V_{1, P^*, \Phi}^*(s_1) - V_{1, P^*, \Phi}^{\hat{\pi}}(s_1) \leq \varepsilon \quad (1)$$

Question:

“Can we design a provably sample-efficient robust RL algorithm using interactive data collection in the training environment?”

Data Generation Mechanism

Interaction protocol:

- ▶ Interact with training Env. P^* for $K \in [N]$ episodes.
- ▶ In each episode $k \in [K]$, use policy π^k to collect data $\{(s_h^k, a_h^k, r_h^k, s_{h+1}^k)\}_{h=1}^H$.
- ▶ After episode k , use historical data to update policy to π^{k+1}

Metric:

- ▶ Online regret:

$$\text{Regret}_{\Phi}(K) := \sum_{k=1}^K V_{1,P^*,\Phi}^*(s_1) - V_{1,P^*,\Phi}^{\pi^k}(s_1).$$

- ▶ Sample complexity: # episodes of interaction suffice to learn ε -optimal robust policy, i.e.

$$V_{1,P^*,\Phi}^*(s_1) - V_{1,P^*,\Phi}^{\hat{\pi}}(s_1) \leq \varepsilon \quad (1)$$

Question:

“Can we design a provably sample-efficient robust RL algorithm using interactive data collection in the training environment?”

Outline

Background and Problem Setup

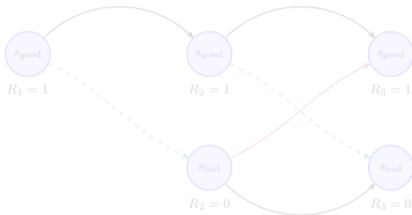
Hardness Result under Interactive Data Collection

Vanishing Minimal Value Assumption and Algorithm Design

Future Works

Hardness Result

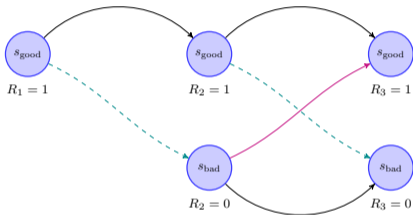
- ▶ For RMDP with total variation (TV) robust set, robust RL with interactive data collection inevitably incurs a regret of $\text{Regret}(K) \geq \Omega(\rho \cdot HK)$, where $\rho =$ size of robust set.
- ▶ Hard example: a simple class of two RMDPs: $\mathcal{S} = \{s_{\text{good}}, s_{\text{bad}}\}$, $\mathcal{A} = \{0, 1\}$.



- Solid lines: transitions of the nominal transition kernel P^* .
 - Dashed lines: transitions of the worst case transition in the robust set.
 - Red solid line: the transition where the two RMDP instances differ in that different action leads to higher transition probability from s_{bad} to s_{good} .
- ▶ When starting from $s_1 = s_{\text{good}}$, the nominal transition kernel keeps the agent at s_{good} and no information at s_{bad} is revealed!

Hardness Result

- ▶ For RMDP with total variation (TV) robust set, robust RL with interactive data collection inevitably incurs a regret of $\text{Regret}(K) \geq \Omega(\rho \cdot HK)$, where $\rho =$ size of robust set.
- ▶ Hard example: a simple class of two RMDPs: $\mathcal{S} = \{s_{\text{good}}, s_{\text{bad}}\}$, $\mathcal{A} = \{0, 1\}$.



- Solid lines: transitions of the nominal transition kernel P^* .
 - Dashed lines: transitions of the worst case transition in the robust set.
 - Red solid line: the transition where the two RMDP instances differ in that different action leads to higher transition probability from s_{bad} to s_{good} .
- ▶ When starting from $s_1 = s_{\text{good}}$, the nominal transition kernel keeps the agent at s_{good} and no information at s_{bad} is revealed!

Outline

Background and Problem Setup

Hardness Result under Interactive Data Collection

Vanishing Minimal Value Assumption and Algorithm Design

Future Works

Combating Hard Instance: Our Assumption

- ▶ We identify the obstacle as the curse of support shift: the disjointedness of the distributional support between training / testing environments.
 - A broader understanding: the states often appearing in testing Env. are extremely hard to arrive in the training Env. Conjecture: imply hardness for other ϕ -divergence robust set.
- ▶ To rule out such instances, we propose the **Vanishing Minimal Value (VMV)** assumption:

Assumption 1 (Vanishing Minimal Value).

The optimal robust value is zero at a specific state, i.e.,

$$\min_{s \in \mathcal{S}} V_{1, P^*, \Phi}^*(s) = 0.$$

- ▶ Rules out the above hard instances (An equivalent form of robust Bellman equation where no explicit support shift occurs)
- ▶ Example (fail state assumption): $\exists s_f \in \mathcal{S}$, s.t. $R_h(s_f, \cdot) = 0$ and $P_h^*(s_f | s_f, \cdot) = 1$.

Combating Hard Instance: Our Assumption

- ▶ We identify the obstacle as the curse of support shift: the disjointedness of the distributional support between training / testing environments.
 - A broader understanding: the states often appearing in testing Env. are extremely hard to arrive in the training Env. Conjecture: imply hardness for other ϕ -divergence robust set.
- ▶ To rule out such instances, we propose the **Vanishing Minimal Value (VMV)** assumption:

Assumption 1 (Vanishing Minimal Value).

The optimal robust value is zero at a specific state, i.e.,

$$\min_{s \in \mathcal{S}} V_{1, P^*, \Phi}^*(s) = 0.$$

- ▶ Rules out the above hard instances (An equivalent form of robust Bellman equation where no explicit support shift occurs)
- ▶ Example (fail state assumption): $\exists s_f \in \mathcal{S}$, s.t. $R_h(s_f, \cdot) = 0$ and $P_h^*(s_f | s_f, \cdot) = 1$.

Algorithm Design

Our algorithm: OPtimistic RObust Value Iteration for TV Robust Set (**OPROVI-TV**).

In each episode $k \in [K]$, it has three stages:

- ▶ (Stage 1: Training env. model estimation) estimate training env. P^* as \hat{P}
- ▶ (Stage 2: Optimistic robust planning) solve the optimal robust policy π^k for the estimated model \hat{P} based on a sophisticated joint consideration of:
 - Robust optimal Bellman equation to ensure exploitation and distributional robustness
 - Optimistic robust value estimation to encourage exploration
- ▶ (Stage 3: Interactive data collection) use policy π^k to collect data.

Algorithm Design

Our algorithm: OPtimistic RObust Value Iteration for TV Robust Set (**OPROVI-TV**).

In each episode $k \in [K]$, it has three stages:

- ▶ (Stage 1: Training env. model estimation) estimate training env. P^* as \hat{P}
- ▶ (Stage 2: Optimistic robust planning) solve the optimal robust policy π^k for the estimated model \hat{P} based on a sophisticated joint consideration of:
 - Robust optimal Bellman equation to ensure exploitation and distributional robustness
 - Optimistic robust value estimation to encourage exploration
- ▶ (Stage 3: Interactive data collection) use policy π^k to collect data.

Algorithm Design

Our algorithm: OPtimistic RObust Value Iteration for TV Robust Set (**OPROVI-TV**).

In each episode $k \in [K]$, it has three stages:

- ▶ (Stage 1: Training env. model estimation) estimate training env. P^* as \hat{P}
- ▶ (Stage 2: Optimistic robust planning) solve the optimal robust policy π^k for the estimated model \hat{P} based on a sophisticated joint consideration of:
 - Robust optimal Bellman equation to ensure exploitation and distributional robustness
 - Optimistic robust value estimation to encourage exploration
- ▶ (Stage 3: Interactive data collection) use policy π^k to collect data.

Algorithm Design

Our algorithm: OPtimistic RObust Value Iteration for TV Robust Set (**OPROVI-TV**).

In each episode $k \in [K]$, it has three stages:

- ▶ (Stage 1: Training env. model estimation) estimate training env. P^* as \hat{P}
- ▶ (Stage 2: Optimistic robust planning) solve the optimal robust policy π^k for the estimated model \hat{P} based on a sophisticated joint consideration of:
 - Robust optimal Bellman equation to ensure exploitation and distributional robustness
 - Optimistic robust value estimation to encourage exploration
- ▶ (Stage 3: Interactive data collection) use policy π^k to collect data.

Theoretical Results

Theorem 1 (Online Regret of OPROVI-TV).

Under the VMV assumption, for any $\rho \in [0, 1)$, OPROVI-TV has an online regret of

$$\text{Regret}(K) \leq \tilde{O}\left(\sqrt{\min\{H, \rho^{-1}\} \cdot H^2 SAK}\right).$$

As a corollary, OPROVI-TV is capable of finding an ε -optimal robust policy within

$$\tilde{O}\left(\min\{H, \rho^{-1}\} \cdot \frac{H^2 SA}{\varepsilon^2}\right)$$

interactive samples. $\tilde{O}(\cdot)$ hides logarithmic factors.

- ▶ First result of this kind
- ▶ Matching the sample complexity lower bound of generative model case [Shi et al., 2023]
- ▶ Requires less sample as the robust set size ρ increases

Theoretical Results

Theorem 1 (Online Regret of OPROVI-TV).

Under the VMV assumption, for any $\rho \in [0, 1)$, OPROVI-TV has an online regret of

$$\text{Regret}(K) \leq \tilde{O}\left(\sqrt{\min\{H, \rho^{-1}\} \cdot H^2 SAK}\right).$$

As a corollary, OPROVI-TV is capable of finding an ε -optimal robust policy within

$$\tilde{O}\left(\min\{H, \rho^{-1}\} \cdot \frac{H^2 SA}{\varepsilon^2}\right)$$

interactive samples. $\tilde{O}(\cdot)$ hides logarithmic factors.

- ▶ First result of this kind
- ▶ Matching the sample complexity lower bound of generative model case [Shi et al., 2023]
- ▶ Requires less sample as the robust set size ρ increases

Outline

Background and Problem Setup

Hardness Result under Interactive Data Collection

Vanishing Minimal Value Assumption and Algorithm Design

Future Works

▶ Function approximations

- Existing works have preliminary results on linear function approximations
- General function approximations: identifying general learnability principle like Bellman rank, bilinear class, Bellman-Eluder dimension, generalized Eluder coefficient for standard online RL.

▶ Other types of robust set

- KL divergence? It is even unknown whether robust RL with interactive data collection is possible in this case.

▶ Robust Markov games

Remark

Provably sample-efficient algorithms exist for generative model/offline learning setup [Blanchet et al., 2023], but remain unknown for interactive data collection!

▶ Function approximations

- Existing works have preliminary results on linear function approximations
- General function approximations: identifying general learnability principle like Bellman rank, bilinear class, Bellman-Eluder dimension, generalized Eluder coefficient for standard online RL.

▶ Other types of robust set

- KL divergence? It is even unknown whether robust RL with interactive data collection is possible in this case.

▶ Robust Markov games

Remark

Provably sample-efficient algorithms exist for generative model/offline learning setup [Blanchet et al., 2023], but remain unknown for interactive data collection!

Thanks for your attention!

References I

- J. Blanchet, M. Lu, T. Zhang, and H. Zhong. Double pessimism is provably efficient for distributionally robust offline reinforcement learning: Generic algorithm and robust partial coverage. Advances in Neural Information Processing Systems, 36, 2023.
- L. Shi, G. Li, Y. Wei, Y. Chen, M. Geist, and Y. Chi. The curious price of distributional robustness in reinforcement learning with a generative model. Advances in Neural Information Processing Systems, 36, 2023.