

# Robust Reinforcement Learning under Model Uncertainty

Shaofeng Zou

University at Buffalo, the State University of New York

6th Workshop on Cognition & Control

Jan. 27th, 2023

- 1 Introduction
- 2 Robust Average-Cost RL
  - Model-based methods
  - Model-free methods
- 3 Summary

- 1 Introduction
- 2 Robust Average-Cost RL
  - Model-based methods
  - Model-free methods
- 3 Summary

# Challenge of Model Mismatch



≠



Training environment  $\neq$  test environment

⇒ Model mismatch

⇒ Severe performance degradation

# Challenge of Model Mismatch



≠



Training environment  $\neq$  test environment

⇒ Model mismatch

⇒ Severe performance degradation

- modeling error between simulator and real-world applications
- non-stationary environment
- unexpected perturbations and potential adversarial attacks

# Challenge of Model Mismatch



≠



Training environment  $\neq$  test environment

⇒ Model mismatch

⇒ Severe performance degradation

- modeling error between simulator and real-world applications
- non-stationary environment
- unexpected perturbations and potential adversarial attacks

## Robust RL:

Find good policy that performs well under model mismatch

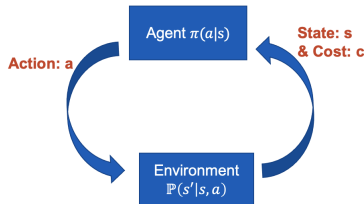
# Markov Decision Processes

- An agent interacts with a stochastic environment: Markov Decision Process (MDP)
- MDP  $(\mathcal{S}, \mathcal{A}, P, c)$

# Markov Decision Processes

- An agent interacts with a stochastic environment: Markov Decision Process (MDP)
- MDP  $(\mathcal{S}, \mathcal{A}, P, c)$

- $\mathcal{S}$ : state space
- $\mathcal{A}$ : action space
- $P$ : transition kernel
- $c$ : cost function



- A stationary policy  $\pi(a|s)$  is a conditional distribution over  $\mathcal{A}$



- Discounted value function for policy  $\pi$  at state  $s$ :

$$\mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) | S_0 = s, \pi \right]$$

- Discounted value function for policy  $\pi$  at state  $s$ :

$$\mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) | S_0 = s, \pi \right]$$

- Goal: find an optimal policy that minimizes value function

$$\min_{\pi} \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) | S_0 = s, \pi \right]$$

- Robust MDP:  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, c, \gamma)$ 
  - $\mathcal{P}$ : uncertainty set of transition kernels
  - Transition kernel at each time step comes from  $\mathcal{P}$ :  
 $\kappa = (P_0, P_1, \dots) \in \otimes_{t \geq 0} \mathcal{P}$

- Robust MDP:  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, c, \gamma)$ 
  - $\mathcal{P}$ : uncertainty set of transition kernels
  - Transition kernel at each time step comes from  $\mathcal{P}$ :  
 $\kappa = (P_0, P_1, \dots) \in \otimes_{t \geq 0} \mathcal{P}$
- Pessimistic approach in face of uncertainty:
  - Worst-case overall cost over uncertainty set
  - Robust value function under the discounted setting:

$$V_{\mathcal{P}, \gamma}^{\pi}(s) = \max_{\kappa \in \otimes_{t \geq 0} \mathcal{P}} \mathbb{E}_{\kappa} \left[ \sum_{t=0}^{\infty} \gamma^t c(S_t, A_t) \mid S_0 = s, \pi \right]$$

- Goal: Optimize the worst-case performance:

$$\min_{\pi} V_{\mathcal{P}, \gamma}^{\pi}(s), \forall s \in \mathcal{S}$$

## Related Works: Robust Discounted RL

- Model-based method:  $\mathcal{P}$  is known
  - e.g., Bagnell et al. (2001); Nilim and El Ghaoui (2004); Iyengar (2005); Wiesemann et al. (2013); Tamar et al. (2014): robust dynamic programming

# Related Works: Robust Discounted RL

- Model-based method:  $\mathcal{P}$  is known
  - e.g., Bagnell et al. (2001); Nilim and El Ghaoui (2004); Iyengar (2005); Wiesemann et al. (2013); Tamar et al. (2014): robust dynamic programming
- Model-free method:  $\mathcal{P}$  is unknown, samples from nominal transition kernel are available
  - Roy et al. (2017); Panaganti et al. (2022): *ellipsoid-structure* uncertainty set
    - convergence needs discount factor bounded away from 1
  - Liu et al. (2022): KL divergence uncertainty set
    - generative model, tabular setting
  - Wang and Zou (2021, 2022) (our work): R-contamination model
    - asymptotic convergence with sample complexity analysis
    - function approximation
  - Zhou et al. (2021); Yang et al. (2021): offline and tabular

$$g_{\mathcal{P}}^{\pi}(s) = \max_{\kappa \in \otimes_{t \geq 0} \mathcal{P}} \lim_{n \rightarrow \infty} \mathbb{E}_{\kappa} \left[ \frac{1}{n} \sum_{t=0}^{n-1} c(S_t, A_t) | S_0 = s, \pi \right]$$

- Model-based method:  $\mathcal{P}$  is known
  - Tewari and Bartlett (2007): bounded-interval uncertainty set, limit method
  - Wang et al. (2023b) (our work): general uncertainty set, robust average-cost Bellman equation, limit method and direct method

$$g_{\mathcal{P}}^{\pi}(s) = \max_{\kappa \in \otimes_{t \geq 0} \mathcal{P}} \lim_{n \rightarrow \infty} \mathbb{E}_{\kappa} \left[ \frac{1}{n} \sum_{t=0}^{n-1} c(S_t, A_t) | S_0 = s, \pi \right]$$

- Model-based method:  $\mathcal{P}$  is known
  - Tewari and Bartlett (2007): bounded-interval uncertainty set, limit method
  - Wang et al. (2023b) (our work): general uncertainty set, robust average-cost Bellman equation, limit method and direct method
- Model-free method:  $\mathcal{P}$  is unknown
  - Wang et al. (2023a) (our work): general uncertainty set, robust relative value iteration with convergence guarantee



- Adversarial **state transition** perturbation: an adversary perturbs the state transition: Vinitzky et al. (2020); Pinto et al. (2017); Abdullah et al. (2019); Hou et al. (2020); Rajeswaran et al. (2017); Atkeson and Morimoto (2003); Morimoto and Doya (2005)

- Adversarial **state transition** perturbation: an adversary perturbs the state transition: Vinitzky et al. (2020); Pinto et al. (2017); Abdullah et al. (2019); Hou et al. (2020); Rajeswaran et al. (2017); Atkeson and Morimoto (2003); Morimoto and Doya (2005)
- Adversarial **sample** perturbation: an adversary modifies state observations Huang et al. (2017); Kos and Song (2017); Lin et al. (2017); Pattanaik et al. (2018); Mandlekar et al. (2017)

- 1 Introduction
- 2 Robust Average-Cost RL
  - Model-based methods
  - Model-free methods
- 3 Summary

# Robust Average-Cost RL

Recall the robust average-cost:

$$g_{\mathcal{P}}^{\pi}(s) \triangleq \lim_{n \rightarrow \infty} \mathbb{E}_{\mathcal{P}} \left[ \frac{1}{n} \sum_{t=0}^{n-1} c_t \mid S_0 = s, \pi \right]$$

$$g_{\mathcal{P}}^{\pi} = \max_{\mathcal{P} \in \mathcal{P}} g_{\mathcal{P}}^{\pi}$$

Goal: Find  $\pi^* = \arg \min_{\pi} g_{\mathcal{P}}^{\pi}$

- Fundamental understanding of robust average-cost MDPs
  - robust average-cost Bellman equation
- Model-based methods:
  - Limit method
  - Direct method
- Model-free methods: robust TD and robust Q-learning

# Robust Average-Cost MDP: Limit Method

## Non-robust setting

(Puterman (1994) point-wise convergence) For any fixed  $P$  and  $\pi$ ,  
 $\lim_{\gamma \rightarrow 1} (1 - \gamma) V_{P, \gamma}^{\pi} = g_P^{\pi}$

- Under non-robust setting, average-cost can be approximated by discounted value function

# Robust Average-Cost MDP: Limit Method

## Non-robust setting

(Puterman (1994) point-wise convergence) For any fixed  $P$  and  $\pi$ ,  
 $\lim_{\gamma \rightarrow 1} (1 - \gamma) V_{P, \gamma}^{\pi} = g_P^{\pi}$

- Under non-robust setting, average-cost can be approximated by discounted value function
- In robust MDP, does it hold that

$$\lim_{\gamma \rightarrow 1} (1 - \gamma) V_{P, \gamma}^{\pi} = g_P^{\pi} \quad ?$$

# Robust Average-Cost MDP: Limit Method

## Non-robust setting

(Puterman (1994) point-wise convergence) For any fixed  $P$  and  $\pi$ ,  
 $\lim_{\gamma \rightarrow 1} (1 - \gamma) V_{P, \gamma}^{\pi} = g_P^{\pi}$

- Under non-robust setting, average-cost can be approximated by discounted value function
- In robust MDP, does it hold that

$$\lim_{\gamma \rightarrow 1} (1 - \gamma) V_{P, \gamma}^{\pi} = g_P^{\pi} \quad ?$$

- Robust discounted Bellman operator (Nilim and El Ghaoui, 2004; Iyengar, 2005):  $\mathbf{T}V = c + \gamma \sum_a \pi(a|s) \sigma_{P_s^a}(V)$ , where  $\sigma_{P_s^a}(V) = \max_{p \in P_s^a} p^{\top} V$  is support function
- $\mathbf{T}$  is a  $\gamma$ -contraction and has  $V_{P, \gamma}^{\pi}$  as its unique fixed point:  $\mathbf{T}V_{P, \gamma}^{\pi} = V_{P, \gamma}^{\pi}$



# Tewari and Bartlett (2007): Bounded-interval Uncertainty Set

- Number of possible worst-case transition kernels is finite
  - Proof of this argument relies on structure of bounded-interval
- Then,  $\min_P$  and  $\lim_{\gamma}$  are interchangeable
- Not generalizable to general uncertainty sets

# Robust Average-Cost MDP: Limit Method

Theorem: uniform convergence

$$\lim_{\gamma \rightarrow 1} (1 - \gamma) V_{P, \gamma}^{\pi} = g_P^{\pi} \text{ uniformly}$$

# Robust Average-Cost MDP: Limit Method

## Theorem: uniform convergence

$$\lim_{\gamma \rightarrow 1} (1 - \gamma) V_{\mathbf{P}, \gamma}^{\pi} = g_{\mathbf{P}}^{\pi} \text{ uniformly}$$

- Then  $\min_{\mathbf{P}}$  and  $\lim_{\gamma \rightarrow 1}$  are interchangeable:

$$\begin{aligned} g_{\mathcal{P}}^{\pi} &= \min_{\mathbf{P} \in \mathcal{P}} g_{\mathbf{P}}^{\pi} \\ &= \min_{\mathbf{P} \in \mathcal{P}} \lim_{\gamma \rightarrow 1} (1 - \gamma) V_{\mathbf{P}, \gamma}^{\pi} = \lim_{\gamma \rightarrow 1} \min_{\mathbf{P} \in \mathcal{P}} (1 - \gamma) V_{\mathbf{P}, \gamma}^{\pi} \\ &= \lim_{\gamma \rightarrow 1} (1 - \gamma) V_{\mathcal{P}, \gamma}^{\pi} \end{aligned}$$

# Robust Average-Cost MDP: Limit Method

## Theorem: uniform convergence

$$\lim_{\gamma \rightarrow 1} (1 - \gamma) V_{\mathbf{P}, \gamma}^{\pi} = g_{\mathbf{P}}^{\pi} \text{ uniformly}$$

- Then  $\min_{\mathbf{P}}$  and  $\lim_{\gamma \rightarrow 1}$  are interchangeable:

$$\begin{aligned} g_{\mathcal{P}}^{\pi} &= \min_{\mathbf{P} \in \mathcal{P}} g_{\mathbf{P}}^{\pi} \\ &= \min_{\mathbf{P} \in \mathcal{P}} \lim_{\gamma \rightarrow 1} (1 - \gamma) V_{\mathbf{P}, \gamma}^{\pi} = \lim_{\gamma \rightarrow 1} \min_{\mathbf{P} \in \mathcal{P}} (1 - \gamma) V_{\mathbf{P}, \gamma}^{\pi} \\ &= \lim_{\gamma \rightarrow 1} (1 - \gamma) V_{\mathcal{P}, \gamma}^{\pi} \end{aligned}$$

- Robust average-cost can be approximated by discounted robust value functions as  $\gamma \rightarrow 1$

# Robust Average-Cost MDP: Limit Method

Basic idea of limit method:

- Set  $\gamma_t \rightarrow 1$
- Apply one-step robust discounted Bellman operator

## Robust value iteration for policy evaluation

**INPUT:**  $\pi, V_0(s) = 0, \forall s, T$

**FOR**  $t = 0, 1, \dots, T - 1$

$$\gamma_t \leftarrow \frac{t+1}{t+2}$$

**FOR** all  $s \in \mathcal{S}$

$$V_{t+1}(s) \leftarrow \mathbb{E}_\pi[(1 - \gamma_t)c(s, A) + \gamma_t \sigma_{\mathcal{P}_s^A}(V_t)]$$

**OUTPUT:**  $V_T$

# Robust Average-Cost MDP: Limit Method

## Convergence of Robust Value Iteration

RVI algorithm converges to robust average-cost:  $\lim_{T \rightarrow \infty} V_T \rightarrow g_{\mathcal{P}}^{\pi}$

# Robust Average-Cost MDP: Limit Method

## Convergence of Robust Value Iteration

RVI algorithm converges to robust average-cost:  $\lim_{T \rightarrow \infty} V_T \rightarrow g_{\mathcal{P}}^{\pi}$

- Solves the policy evaluation problem under the robust average-cost setting

## Convergence of Robust Value Iteration

RVI algorithm converges to robust average-cost:  $\lim_{T \rightarrow \infty} V_T \rightarrow g_{\mathcal{P}}^{\pi}$

- Solves the policy evaluation problem under the robust average-cost setting
- Convergence rate:  $\|V_T - g_{\mathcal{P}}^{\pi}\| = \mathcal{O}\left(\frac{1}{T}\right)$



# Robust Average-Cost MDP: Limit Method

- The limit method also works for optimal control problems

## Robust value iteration for optimal control

**INPUT:**  $V_0(s) = 0, \forall s, T$

**FOR**  $t = 0, 1, \dots, T - 1$

$$\gamma_t \leftarrow \frac{t+1}{t+2}$$

**FOR** all  $s \in \mathcal{S}$

$$V_{t+1}(s) \leftarrow \min_{a \in \mathcal{A}} \{ (1 - \gamma_t)c(s, a) + \gamma_t \sigma_{\mathcal{P}_s^a}(V_t) \}$$

**FOR**  $s \in \mathcal{S}$

$$\pi_t(s) \leftarrow \arg \min_{a \in \mathcal{A}} \{ (1 - \gamma_t)c(s, a) + \gamma_t \sigma_{\mathcal{P}_s^a}(V_t) \}$$

**OUTPUT:**  $\pi_T, V_T$

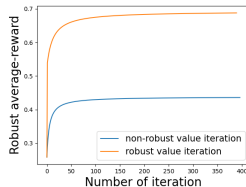
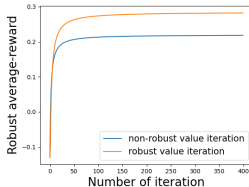
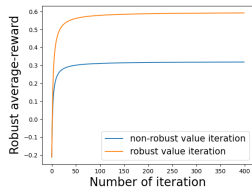
## Convergence of robust value iteration

$$V_T \rightarrow g_{\mathcal{P}}^*, \pi_T \rightarrow \pi^* = \arg \min_{\pi} g_{\mathcal{P}}^{\pi}$$

# Robust Average-Cost MDP: Limit Method

Non-robust value iteration vs robust value iteration:

- under different uncertainty sets (contamination model, total variation model and KL-divergence model)
- evaluate **worst-case** performance of obtained policy



RVI is more robust than non-robust value iteration method

# Robust Average-Cost MDP: Limit Method

## Robust Blackwell optimality

There exists  $\delta < 1$ , such that for any  $\delta < \gamma < 1$ , if  $\pi^* = \arg \min_{\pi} V_{\mathcal{P}, \gamma}^{\pi}$  is optimal to robust discounted value function, then  $\pi^*$  is also optimal to robust average-cost  $\pi^* \in \arg \min_{\pi} g_{\mathcal{P}}^{\pi}$ .

## Robust Blackwell optimality

There exists  $\delta < 1$ , such that for any  $\delta < \gamma < 1$ , if  $\pi^* = \arg \min_{\pi} V_{\mathcal{P}, \gamma}^{\pi}$  is optimal to robust discounted value function, then  $\pi^*$  is also optimal to robust average-cost  $\pi^* \in \arg \min_{\pi} g_{\mathcal{P}}^{\pi}$ .

- Fundamental relationship between the robust discounted MDPs and robust average-cost MDPs

## Robust Blackwell optimality

There exists  $\delta < 1$ , such that for any  $\delta < \gamma < 1$ , if  $\pi^* = \arg \min_{\pi} V_{\mathcal{P}, \gamma}^{\pi}$  is optimal to robust discounted value function, then  $\pi^*$  is also optimal to robust average-cost  $\pi^* \in \arg \min_{\pi} g_{\mathcal{P}}^{\pi}$ .

- Fundamental relationship between the robust discounted MDPs and robust average-cost MDPs
- Analog to Blackwell optimality of non-robust setting

# Robust Average-Cost MDP: Limit Method

## Robust Blackwell optimality

There exists  $\delta < 1$ , such that for any  $\delta < \gamma < 1$ , if  $\pi^* = \arg \min_{\pi} V_{\mathcal{P},\gamma}^{\pi}$  is optimal to robust discounted value function, then  $\pi^*$  is also optimal to robust average-cost  $\pi^* \in \arg \min_{\pi} g_{\mathcal{P}}^{\pi}$ .

- Fundamental relationship between the robust discounted MDPs and robust average-cost MDPs
- Analog to Blackwell optimality of non-robust setting
- Proofs of non-robust setting and bounded-interval uncertainty set  
Tewari and Bartlett (2007): for two policies  $\pi$  and  $\nu$ :  
 $f_{\pi,\nu}(\gamma) \triangleq V_{\mathcal{P},\gamma}^{\pi} - V_{\mathcal{P},\gamma}^{\nu}$  is rational function, thus has finite many zeros.  
This does not hold in robust setting as  $V_{\mathcal{P},\gamma}^{\pi} - V_{\mathcal{P},\gamma}^{\nu}$  is not rational due to max

# Robust Average-Cost MDP: Limit Method

## The limit method

- solves robust average-cost MDPs using robust discounted MDPs as intermediate steps
- based on robust discounted MDPs, does not directly study the fundamental structure of robust average-cost MDPs

# Robust Average-Cost MDP: Direct Method

## Assumption

The Markov chain induced by any  $P \in \mathcal{P}$  and any  $\pi$  is a unichain.

## Optimal robust Bellman equation

If  $(g, V)$  is a solution to

$$V(s) = \min_a \{c(s, a) - g + \sigma_{\mathcal{P}_s^a}(V)\}, \forall s,$$

then  $g = g_{\mathcal{P}}^*$ . If we further set

$$\pi^*(s) = \arg \min_a \{c(s, a) + \sigma_{\mathcal{P}_s^a}(V)\}$$

for any  $s \in \mathcal{S}$ , then  $\pi^*$  is an optimal robust policy.

- Solving robust average-cost MDPs can be done by solving the robust Bellman equation



# Robust Average-Cost MDP: Direct Method

$$V(s) = \mathbf{T}(V) = \min_a \{c(s, a) - g + \sigma_{\mathcal{P}_s^a}(V)\}$$

How to solve the robust Bellman equation?

# Robust Average-Cost MDP: Direct Method

$$V(s) = \mathbf{T}(V) = \min_a \{c(s, a) - g + \sigma_{\mathcal{P}_s^a}(V)\}$$

How to solve the robust Bellman equation?

- Discounted setting: apply the robust Bellman operator recursively ( $\gamma$ -contraction)

# Robust Average-Cost MDP: Direct Method

$$V(s) = \mathbf{T}(V) = \min_a \{c(s, a) - g + \sigma_{\mathcal{P}_s^a}(V)\}$$

How to solve the robust Bellman equation?

- Discounted setting: apply the robust Bellman operator recursively ( $\gamma$ -contraction)
- **Average setting: not a contraction, may have multiple fixed points and algorithm may diverge**

# Robust Average-Cost MDP: Direct Method

$$V(s) = \mathbf{T}(V) = \min_a \{c(s, a) - g + \sigma_{\mathcal{P}_s^a}(V)\}$$

How to solve the robust Bellman equation?

- Discounted setting: apply the robust Bellman operator recursively ( $\gamma$ -contraction)
- **Average setting: not a contraction, may have multiple fixed points and algorithm may diverge**
- Robust relative value iteration (RRVI)
  - subtract an offset function to keep iterates stable
  - prove it is multi-step contraction

# Robust Average-Cost MDP: Direct Method

## Robust relative value iteration

**INPUT:**  $V_0$ ,  $\epsilon$  and arbitrary  $s^* \in \mathcal{S}$

**WHILE** TRUE

**FOR** all  $s \in \mathcal{S}$

$$V_{t+1}(s) \leftarrow \min_a (c(s, a) + \sigma_{\mathcal{P}_s^a}(V_t) - f(V_t))$$

**OUTPUT:**  $f(V_t), V_t$

- For example  $f(V) = V(s^*)$  for some reference state  $s^*$  and  $f(V)$  is the mean of  $V$ , to "offset" the increase of  $V$

# Robust Average-Cost MDP: Direct Method

Convergence of robust relative value iteration

$(f(V_t), V_t)$  converges to a solution to the optimal robust Bellman equation

## Convergence of robust relative value iteration

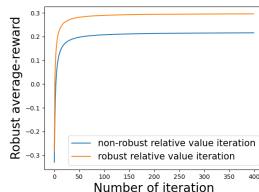
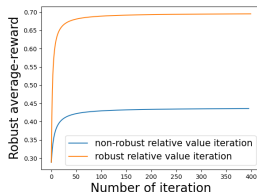
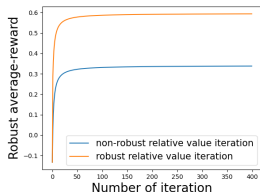
$(f(V_t), V_t)$  converges to a solution to the optimal robust Bellman equation

- Finds a solution to optimal robust Bellman equation and hence optimal robust average-cost and optimal policy
- Linear convergence rate, faster than limit method

# Robust Average-Cost MDP: Limit Method

Non-robust value iteration vs robust value iteration:

- under different uncertainty sets (contamination model, total variation model and KL-divergence model)
- evaluate **worst-case** performance of obtained policy



RRVI is more robust than non-robust relative value iteration



# Model-free Method for Robust Average-Cost RL

- Idea: generalize RVI Q-learning to robust setting
- Major challenges:
  - The Bellman operator for robust average-cost MDPs is not contraction: possible multiple fixed point
  - Construction of unbiased estimator for robust Bellman operator for various uncertainty sets

# Robust Q-Learning for Average RL

## Optimal robust Bellman equation

If  $(g, Q)$  is a solution to the optimal robust Bellman equation

$$Q(s, a) = r(s, a) - g + \sigma_{\mathcal{P}_s^a}(V_Q), \forall s, a,$$

then 1)  $g = g_{\mathcal{P}}^*$ ;

2) the greedy policy w.r.t.  $Q$ :  $\pi_Q(s) = \arg \min_a Q(s, a)$  is an optimal robust policy;

3)  $V_Q(s) \triangleq \min_a Q(s, a) = V_{\mathcal{P}}^{\pi_Q}(s) + ce$  for some  $\mathcal{P} \in \Omega_g^{\pi_Q}$ ,  $c \in \mathbb{R}$ .

- worst-case transition kernels:  $\Omega_g^{\pi} = \{\mathcal{P} \in \mathcal{P} : g_{\mathcal{P}}^{\pi} = g_{\mathcal{P}}^{\pi}\}$
- relative value function:  $V_{\mathcal{P}}^{\pi} = \mathbb{E}_{\pi, \mathcal{P}} \left[ \sum_{t=0}^{\infty} \mathcal{P}^t(r - g_{\mathcal{P}}^{\pi}) \right]$

# Robust Q-Learning for Average RL

## Optimal robust Bellman equation

If  $(g, Q)$  is a solution to the optimal robust Bellman equation

$$Q(s, a) = r(s, a) - g + \sigma_{\mathcal{P}_s^a}(V_Q), \forall s, a,$$

then 1)  $g = g_{\mathcal{P}}^*$ ;

2) the greedy policy w.r.t.  $Q$ :  $\pi_Q(s) = \arg \min_a Q(s, a)$  is an optimal robust policy;

3)  $V_Q(s) \triangleq \min_a Q(s, a) = V_{\mathcal{P}}^{\pi_Q}(s) + ce$  for some  $\mathcal{P} \in \Omega_g^{\pi_Q}$ ,  $c \in \mathbb{R}$ .

- worst-case transition kernels:  $\Omega_g^{\pi} = \{\mathcal{P} \in \mathcal{P} : g_{\mathcal{P}}^{\pi} = g_{\mathcal{P}}^{\pi}\}$
- relative value function:  $V_{\mathcal{P}}^{\pi} = \mathbb{E}_{\pi, \mathcal{P}} \left[ \sum_{t=0}^{\infty} \mathcal{P}^t(r - g_{\mathcal{P}}^{\pi}) \right]$
- non-robust setting: Bellman equation is linear, and thus structure of solutions can be easily characterized
- robust Bellman equation is non-linear

# Robust Q-Learning for Average RL

## Robust RVI Q-learning

**INPUT:**  $Q_0, \alpha_n, N$

**FOR**  $n = 0, \dots, N - 1$

$$Q_{n+1} \leftarrow Q_n + \alpha_n (\hat{\mathbf{H}} Q_n - f(Q_n) - Q_n)$$

**OUTPUT:**  $Q_N$

- $\hat{\mathbf{H}}Q$ : unbiased estimator of  $\mathbf{H}Q = c(s, a) + \sigma_{\mathcal{P}_s^a}(V_Q)$
- $f(Q) : \mathbb{R}^{|\mathcal{SA}|} \rightarrow \mathbb{R}$ : "offset" increase of  $Q_n$  and keep iterates stable
- $f(Q_n)$ : estimator of average-cost  $g_{\mathcal{P}}^*$

# Robust Q-Learning for Average RL

## Convergence of Robust Q-Learning

If  $\hat{\mathbf{H}}$  is unbiased and has bounded variance, then almost surely,

- 1)  $f(Q_n)$  converges to  $g_{\mathcal{P}}^*$ ;
- 2) greedy policy  $\pi_{Q_n}(s) \triangleq \arg \max_a Q_n(s, a)$  converges to an optimal robust average-cost policy.

- To show convergence, we need structure of solution to robust average-cost Bellman equation to characterize the equilibrium of associated ODE, and prove it is globally asymptotically stable

# Robust Q-Learning for Robust Average-Cost RL

- How to construct  $\hat{\mathbf{H}}$ ?

- R-contamination model: MLE method

$$\mathcal{P}_s^a = \{(1 - R)p_s^a + Rq | q \in \Delta_{\mathcal{S}}\}, \text{ for some } 0 \leq R \leq 1$$

- Other uncertainty models, e.g., total variation, Chi-square, Wasserstein distance?
- The support function  $\sigma_{\mathcal{P}}(V)$  w.r.t. general uncertainty sets is **non-linear** in nominal kernel
- MLE method  $\Rightarrow$  biased estimator

# Robust Q-Learning for Average RL

- Multi-level Monte-Carlo method (Blanchet and Glynn, 2015)

For any  $s, a$ :

- Generate  $N$  according to **Geo**( $\Psi$ )
- Sample  $2^{N+1}$  samples:  $\{s'_i\}, i = 1, \dots, 2^{N+1}$
- divide these  $2^{N+1}$  samples into two groups: samples with odd indices, and samples with even indices
- individually calculate the empirical distribution of  $s'$  using the even-index samples, odd-index ones, all the samples, and the first sample:  $\hat{P}_{s,N+1}^{a,E} = \frac{1}{2^N} \sum_{i=1}^{2^N} \mathbb{1}_{s'_{2i}}, \quad \hat{P}_{s,N+1}^{a,O} = \frac{1}{2^N} \sum_{i=1}^{2^N} \mathbb{1}_{s'_{2i-1}}, \quad \hat{P}_{s,N+1}^a = \frac{1}{2^{N+1}} \sum_{i=1}^{2^{N+1}} \mathbb{1}_{s'_i}, \quad \hat{P}_{s,N+1}^{a,1} = \mathbb{1}_{s'_1}$
- Use these estimated transition kernels as nominal kernels to construct four estimated uncertainty sets  $\hat{P}_{s,N+1}^{a,E}, \hat{P}_{s,N+1}^{a,O}, \hat{P}_{s,N+1}^a, \hat{P}_{s,N+1}^{a,1}$

The multi-level estimator is then defined as

$$\hat{\sigma}_{\mathcal{P}_s^a}(V) \triangleq \sigma_{\hat{\mathcal{P}}_{s,N+1}^{a,1}}(V) + \frac{\Delta_N(V)}{p_N}, \quad (1)$$

where  $p_N = \Psi(1 - \Psi)^N$  and

$$\Delta_N(V) \triangleq \sigma_{\hat{\mathcal{P}}_{s,N+1}^a}(V) - \frac{\sigma_{\hat{\mathcal{P}}_{s,N+1}^{a,E}}(V) + \sigma_{\hat{\mathcal{P}}_{s,N+1}^{a,O}}(V)}{2}.$$



For uncertainty sets including contamination model, total variation, Chi-squared divergence, Kullback-Leibler (KL) divergence and Wasserstein distance:

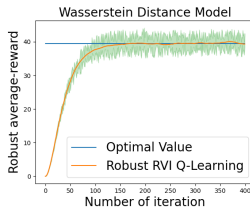
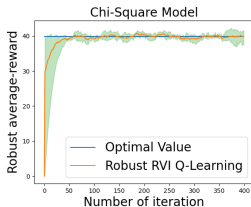
$$\begin{aligned}\mathbb{E}[\hat{\sigma}_{\mathcal{P}_s^a}(V)] &= \sigma_{\mathcal{P}_s^a}(V), \\ \text{Var}[\hat{\sigma}_{\mathcal{P}_s^a}(V)(s)] &\leq C(1 + \|V\|^2).\end{aligned}$$

- For five uncertainty sets above,  $\hat{\sigma}_{\mathcal{P}_s^a}(V)$  is unbiased and has bounded variance
- Implies convergence of robust RVI Q-learning
- Can also be applied to robust discounted setting

# Experiments

## Convergence of robust Q-learning

- Different uncertainty sets, e.g., Chi-Square Model and Wasserstein distance model
- Plot  $f(Q_t)$  (estimate of average reward)
- Baseline is computed using model-based RVI method discussed before



Robust Q-learning converges to the optimal robust average-reward

- 1 Introduction
- 2 Robust Average-Cost RL
  - Model-based methods
  - Model-free methods
- 3 Summary

# Summary

- Robust average-cost RL
- Fundamental understanding
  - Robust average-cost Bellman equation
  - Solution characterization
  - Blackwell optimality
- Model-based approach
  - Limit method
  - Direct method
- Model-free approach: robust RVI Q-learning with convergence guarantee

# Reference I

- Abdullah, M. A., Ren, H., Ammar, H. B., Milenkovic, V., Luo, R., Zhang, M., and Wang, J. (2019). Wasserstein robust reinforcement learning. *arXiv preprint arXiv:1907.13196*.
- Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. (2021). On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):1–76.
- Asadi, K. and Littman, M. L. (2017). An alternative softmax operator for reinforcement learning. In *Proc. International Conference on Machine Learning (ICML)*, volume 70, pages 243–252. JMLR.
- Atkeson, C. G. and Morimoto, J. (2003). Nonparametric representation of policies and value functions: A trajectory-based approach. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, pages 1643–1650.
- Bagnell, J. A., Ng, A. Y., and Schneider, J. G. (2001). Solving uncertain Markov decision processes.

## Reference II

- Bhandari, J. and Russo, D. (2021). On the linear convergence of policy gradient methods for finite MDPs. In *Proc. International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 2386–2394. PMLR.
- Blanchet, J. H. and Glynn, P. W. (2015). Unbiased monte carlo for optimization and functions of expectations via multi-level randomization. In *2015 Winter Simulation Conference (WSC)*, pages 3656–3667. IEEE.
- Bolte, J., Daniilidis, A., and Lewis, A. (2007). The Łojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems. *SIAM Journal on Optimization*, 17(4):1205–1223.
- Cen, S., Cheng, C., Chen, Y., Wei, Y., and Chi, Y. (2021). Fast global convergence of natural policy gradient methods with entropy regularization. *Operations Research*.
- Hou, L., Pang, L., Hong, X., Lan, Y., Ma, Z., and Yin, D. (2020). Robust reinforcement learning with Wasserstein constraint. *arXiv preprint arXiv:2006.00945*.

## Reference III

- Huang, S., Papernot, N., Goodfellow, I., Duan, Y., and Abbeel, P. (2017). Adversarial attacks on neural network policies. In *Proc. International Conference on Learning Representations (ICLR)*.
- Iyengar, G. N. (2005). Robust dynamic programming. *Mathematics of Operations Research*, 30(2):257–280.
- Karimi, H., Nutini, J., and Schmidt, M. (2016). Linear convergence of gradient and proximal-gradient methods under the Polyak-Łojasiewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 795–811. Springer.
- Kos, J. and Song, D. (2017). Delving into adversarial attacks on deep policies. In *Proc. International Conference on Learning Representations (ICLR)*.
- Kruger, A. Y. (2003). On Fréchet subdifferentials. *Journal of Mathematical Sciences*, 116(3):3325–3358.

- Lin, Y.-C., Hong, Z.-W., Liao, Y.-H., Shih, M.-L., Liu, M.-Y., and Sun, M. (2017). Tactics of adversarial attack on deep reinforcement learning agents. In *Proc. International Joint Conferences on Artificial Intelligence (IJCAI)*, pages 3756–3762.
- Liu, Z., Bai, Q., Blanchet, J., Dong, P., Xu, W., Zhou, Z., and Zhou, Z. (2022). Distributionally robust  $q$ -learning. In *International Conference on Machine Learning*, pages 13623–13643. PMLR.
- Mandlekar, A., Zhu, Y., Garg, A., Fei-Fei, L., and Savarese, S. (2017). Adversarially robust policy learning: Active construction of physically-plausible perturbations. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3932–3939. IEEE.
- Morimoto, J. and Doya, K. (2005). Robust reinforcement learning. *Neural computation*, 17(2):335–359.



- Nilim, A. and El Ghaoui, L. (2004). Robustness in Markov decision problems with uncertain transition matrices. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, pages 839–846.
- Panaganti, K., Xu, Z., Kalathil, D., and Ghavamzadeh, M. (2022). Robust reinforcement learning using offline data. *arXiv preprint arXiv:2208.05129*.
- Pattanaik, A., Tang, Z., Liu, S., Bommannan, G., and Chowdhary, G. (2018). Robust deep reinforcement learning with adversarial attacks. In *Proc. International Conference on Autonomous Agents and MultiAgent Systems*, pages 2040–2042.
- Pinto, L., Davidson, J., Sukthankar, R., and Gupta, A. (2017). Robust adversarial reinforcement learning. In *Proc. International Conference on Machine Learning (ICML)*, pages 2817–2826. PMLR.
- Puterman, M. L. (1994). Markov decision processes: Discrete stochastic dynamic programming.

## Reference VI

- Rajeswaran, A., Ghotra, S., Ravindran, B., and Levine, S. (2017). Epopt: Learning robust neural network policies using model ensembles. In *Proc. International Conference on Learning Representations (ICLR)*.
- Roy, A., Xu, H., and Pokutta, S. (2017). Reinforcement learning under model mismatch. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, pages 3046–3055.
- Tamar, A., Mannor, S., and Xu, H. (2014). Scaling up robust MDPs using function approximation. In *Proc. International Conference on Machine Learning (ICML)*, pages 181–189. PMLR.
- Tewari, A. and Bartlett, P. L. (2007). Bounded parameter markov decision processes with average reward criterion. In *International Conference on Computational Learning Theory*, pages 263–277. Springer.
- Vinitisky, E., Du, Y., Parvate, K., Jang, K., Abbeel, P., and Bayen, A. (2020). Robust reinforcement learning using adversarial populations. *arXiv preprint arXiv:2008.01825*.

## Reference VII

- Wang, Y., Velasquez, A., Atia, G., Prater-Bennette, A., and Zou, S. (2023a). Model-free robust average-reward reinforcement learning. In *submitted to ICML*.
- Wang, Y., Velasquez, A., Atia, G., Prater-Bennette, A., and Zou, S. (2023b). Robust average-reward Markov decision processes. In *Proc. Conference on Artificial Intelligence (AAAI)*.
- Wang, Y. and Zou, S. (2021). Online robust reinforcement learning with model uncertainty. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*.
- Wang, Y. and Zou, S. (2022). Policy gradient method for robust reinforcement learning. In *Proc. International Conference on Machine Learning (ICML)*, volume 162, pages 23484–23526. PMLR.
- Wiesemann, W., Kuhn, D., and Rustem, B. (2013). Robust Markov decision processes. *Mathematics of Operations Research*, 38(1):153–183.

- Yang, W., Zhang, L., and Zhang, Z. (2021). Towards theoretical understandings of robust Markov decision processes: Sample complexity and asymptotics. *arXiv preprint arXiv:2105.03863*.
- Zhou, Z., Bai, Q., Zhou, Z., Qiu, L., Blanchet, J., and Glynn, P. (2021). Finite-sample regret bound for distributionally robust offline tabular reinforcement learning. In *Proc. International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 3331–3339. PMLR.

4 Robust Discounted RL

5 Robust Sub-gradient

- Policy evaluation: robust TD (tabular), robust TDC (with function approximation)
- Optimal control: robust Q-learning (value-based), robust policy gradient (policy-based)

# Value-Based Optimal Control: Robust Q-Learning

- Goal: Find a policy optimizing the worst-case performance

$$Q_{\mathcal{P},\gamma}^{\pi}(s, a) = \max_{\kappa \in \otimes_{t \geq 0} \mathcal{P}} \mathbb{E}_{\kappa} \left[ \sum_{t=0}^{\infty} \gamma^t c(S_t, A_t) | S_0 = s, A_0 = a, \pi \right]$$

$$Q_{\mathcal{P},\gamma}^*(s, a) = \min_{\pi} \max_{\kappa \in \otimes_{t \geq 0} \mathcal{P}} \mathbb{E}_{\kappa} \left[ \sum_{t=0}^{\infty} \gamma^t c(S_t, A_t) | S_0 = s, A_0 = a, \pi \right]$$

- Finding  $\pi^*$  is equivalent to find  $Q_{\mathcal{P},\gamma}^*$

# Value-Based Optimal Control: Robust Q-Learning

- **Robust Bellman operator** (Nilim and El Ghaoui, 2004):  
 $\mathbf{T}Q(s, a) = c(s, a) + \gamma \sigma_{\mathcal{P}_s^a}(\min_{a \in \mathcal{A}} Q(s, a))$ , where  
 $\sigma_{\mathcal{P}}(v) = \max_{p \in \mathcal{P}} p^\top v$
- $\mathbf{T}$  is a  $\gamma$ -contraction and has  $Q_{\mathcal{P}, \gamma}^*$  as its unique fixed point:  $\mathbf{T}Q^* = Q^*$
- Idea: recursively apply  $\mathbf{T}$



# Value-Based Optimal Control: Robust Q-Learning

- **Robust Bellman operator** (Nilim and El Ghaoui, 2004):  
$$\mathbf{T}Q(s, a) = c(s, a) + \gamma \sigma_{\mathcal{P}_s^a}(\min_{a \in \mathcal{A}} Q(s, a))$$
- Idea: recursively apply  $\mathbf{T}$
- Model-free setting:
  - No information about the environment or the uncertainty set  $\mathcal{P}$
  - Samples are generated under the nominal environment, generally is different from the worst-case environment
- Estimated the support function  $\sigma_{\mathcal{P}_s^a}(Q)$  using the nominal samples

# R-Contamination Uncertainty Sets

In this work, we mainly focus on  $R$ -contamination uncertainty set:

- $\mathcal{P}_s^a = \{(1 - R)p_s^a + Rq | q \in \Delta_{|\mathcal{S}|}\}, s \in \mathcal{S}, a \in \mathcal{A}$ , for some  $0 \leq R \leq 1$
- Adversarial model: nature can arbitrarily modify transition kernel with probability  $R$

For nominal sample  $O_t = (s_t, a_t, s_{t+1})$ :

- Maximum likelihood estimation of transition kernel  $\hat{p}_t \triangleq \mathbb{1}_{s_{t+1}}$
- Estimated uncertainty set  $\hat{\mathcal{P}}_t \triangleq \{(1 - R)\hat{p}_t + Rq | q \in \Delta_{|S|}\}$
- Compute the support function w.r.t.  $\hat{\mathcal{P}}_t$ :  
$$\sigma_{\hat{\mathcal{P}}_t}(V_t) = (1 - R)V_t(s_{t+1}) + R \max_s V_t(s)$$
- Update Q-function  
$$Q_{t+1}(s_t, a_t) \leftarrow (1 - \alpha_t)Q_t(s_t, a_t) + \alpha_t(c_t + \gamma \sigma_{\hat{\mathcal{P}}_t}(\min_a Q_t))$$

# Robust Q-learning

**Initialization:**  $T$ ,  $Q_0(s, a)$  for all  $(s, a)$ , behavior policy  $\pi_b$ ,  $s_0$ , step size  $\alpha_t$

**For**  $t = 0, 1, 2, \dots, T - 1$

    Choose  $a_t$  according to  $\pi_b(\cdot|s_t)$

    Observe  $s_{t+1}$  and  $c_t$

$V_t(s) \leftarrow \min_{a \in \mathcal{A}} Q_t(s, a)$ ,  $\forall s \in \mathcal{S}$

$Q_{t+1}(s_t, a_t) \leftarrow (1 - \alpha_t)Q_t(s_t, a_t) + \alpha_t(c_t + \gamma \sigma_{\hat{p}_t}(V_t))$

$Q_{t+1}(s, a) \leftarrow Q_t(s, a)$  for  $(s, a) \neq (s_t, a_t)$

**Output:**  $Q_T$

## Theorem

*(Asymptotic Convergence) If step sizes  $\alpha_t$  satisfy that  $\sum_{t=0}^{\infty} \alpha_t = \infty$  and  $\sum_{t=0}^{\infty} \alpha_t^2 < \infty$ , then  $Q_t \rightarrow Q_{\mathcal{P}, \gamma}^*$  as  $t \rightarrow \infty$  almost surely.*

# Theoretical Results

Assumption: The Markov chain induced by behavior policy  $\pi_b$  and transition kernel  $p_s^a$  is uniformly ergodic

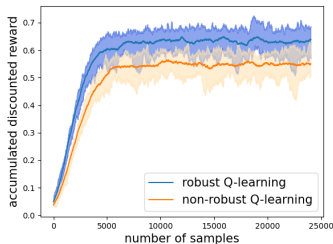
## Theorem

(Finite-Time Error Bound) For any  $\epsilon$ , set  $T = \tilde{\mathcal{O}}\left(\frac{1}{\mu_{\min}(1-\gamma)^5\epsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}\right)$ , then  $\|Q_T - Q_{\mathcal{P},\gamma}^*\| \leq \mathcal{O}(\epsilon)$ .

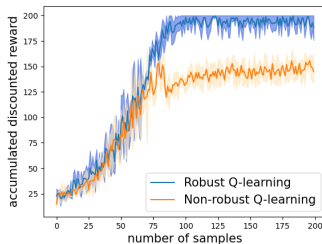
- Matches the sample complexity of non-robust Q-learning (up to some constants)
- First online, model-free method for robust RL with sample complexity result

# Experiments on Robust Q-Learning

Train Q-learning and robust Q-learning under a uniformly perturbed MDP  
Test their outputs in the real unperturbed environment  
Robust Q-learning achieves higher reward than Q-learning



(a) FrozenLake



(b) Cartpole

# Summary on Robust Q-Learning

- For R-contamination model, use maximum likelihood estimation as the estimated nominal transition kernel, and define the estimated uncertainty set
- The support function w.r.t. the estimated uncertainty set is unbiased
- This method can be also applied to policy evaluation problem, e.g., robust TD (tabular case) or robust TDC (function approximation case)



Value-based method:

- Obtains the optimal policy using the robust value functions as an intermediate step, not direct
- Has great memory cost when the problem scale is large

**Our work:** Direct policy search method with global optimality for model-free robust RL problems, and further characterize its sample complexity

**Robust value function  $V_{\mathcal{P},\gamma}^\pi$  may not be differentiable and non-convex**  
$$V_{\mathcal{P},\gamma}^\pi(s) = \max_{\kappa \in \otimes_{t \geq 0} \mathcal{P}} \mathbb{E}_\kappa \left[ \sum_{t=0}^{\infty} \gamma^t c(S_t, A_t) \mid S_0 = s, \pi \right]$$

- Generalize the vanilla policy gradient to the robust policy sub-gradient method, which shows global optimality

**In model-free setting, robust value functions measure the worst-case performance and are impossible to estimate using Monte Carlo method**

- Propose a robust TD algorithm (which can be applied together with function approximation) to estimate the value functions, and further develop a robust actor-critic algorithm

# Main Contributions

- Derivation of robust policy gradient:  $\partial V_{\mathcal{P}, \gamma}^{\pi_{\theta}}(s)$
- Global optimality guarantee and finite-time complexity bound
- **Model-free** robust actor-critic, its convergence and sample complexity

- Idea: derive gradient of  $J_\rho(\pi) \triangleq \mathbb{E}_\rho[V_{\mathcal{P},\gamma}^\pi(S)]$ , and run gradient descent

# Robust Policy Gradient

- Idea: derive gradient of  $J_\rho(\pi) \triangleq \mathbb{E}_\rho[V_{\mathcal{P},\gamma}^\pi(S)]$ , and run gradient descent
- Robust value function  $V_{\mathcal{P},\gamma}^\pi$  is not differentiable everywhere because of max over  $\kappa$

$$V_{\mathcal{P},\gamma}^\pi(s) = \max_{\kappa} \mathbb{E}_\kappa \left[ \sum_{t=0}^{\infty} \gamma^t c(S_t, A_t) | S_0 = s, \pi \right]$$

# Robust Policy Gradient

- Idea: derive gradient of  $J_\rho(\pi) \triangleq \mathbb{E}_\rho[V_{\mathcal{P},\gamma}^\pi(S)]$ , and run gradient descent
- Robust value function  $V_{\mathcal{P},\gamma}^\pi$  is not differentiable everywhere because of max over  $\kappa$

$$V_{\mathcal{P},\gamma}^\pi(s) = \max_{\kappa} \mathbb{E}_\kappa \left[ \sum_{t=0}^{\infty} \gamma^t c(S_t, A_t) | S_0 = s, \pi \right]$$

- Major challenge lies in the max operator

# Robust Policy Sub-gradient

- Consider a parametric policy class  $\Pi_{\Theta} = \{\pi_{\theta} : \theta \in \Theta\}$

## Theorem (Robust Policy Sub-gradient)

Define

$$\begin{aligned}\psi_{\rho}(\theta) \triangleq & \frac{\gamma R}{(1-\gamma)(1-\gamma+\gamma R)} \sum_{s \in \mathcal{S}} d_{s\theta}^{\pi_{\theta}}(s) \sum_{a \in \mathcal{A}} \nabla \pi_{\theta}(a|s) Q_{\mathcal{P},\gamma}^{\pi_{\theta}}(s, a) \\ & + \frac{1}{1-\gamma+\gamma R} \sum_{s \in \mathcal{S}} d_{\rho}^{\pi_{\theta}}(s) \sum_{a \in \mathcal{A}} \nabla \pi_{\theta}(a|s) Q_{\mathcal{P},\gamma}^{\pi_{\theta}}(s, a),\end{aligned}$$

then (1) almost everywhere in  $\Theta$ ,  $J_{\rho}(\theta)$  is differentiable and

$\psi_{\rho}(\theta) = \nabla J_{\rho}(\theta)$ ;

(2) at non-differentiable  $\theta$ ,  $\psi_{\rho}(\theta) \in \partial J_{\rho}(\theta)$ .

- $\partial J_{\rho}(\theta)$ : set of Fréchet sub-differential (Kruger, 2003) of  $J_{\rho}$  at  $\theta$



# Robust Policy Sub-gradient

- Consider a parametric policy class  $\Pi_{\Theta} = \{\pi_{\theta} : \theta \in \Theta\}$

## Theorem (Robust Policy Sub-gradient)

*Define*

$$\begin{aligned}\psi_{\rho}(\theta) \triangleq & \frac{\gamma R}{(1-\gamma)(1-\gamma+\gamma R)} \sum_{s \in \mathcal{S}} d_{s_{\theta}}^{\pi_{\theta}}(s) \sum_{a \in \mathcal{A}} \nabla \pi_{\theta}(a|s) Q_{\mathcal{P}, \gamma}^{\pi_{\theta}}(s, a) \\ & + \frac{1}{1-\gamma+\gamma R} \sum_{s \in \mathcal{S}} d_{\rho}^{\pi_{\theta}}(s) \sum_{a \in \mathcal{A}} \nabla \pi_{\theta}(a|s) Q_{\mathcal{P}, \gamma}^{\pi_{\theta}}(s, a),\end{aligned}$$

*then (1) almost everywhere in  $\Theta$ ,  $J_{\rho}(\theta)$  is differentiable and*

*$\psi_{\rho}(\theta) = \nabla J_{\rho}(\theta)$ ;*

*(2) at non-differentiable  $\theta$ ,  $\psi_{\rho}(\theta) \in \partial J_{\rho}(\theta)$ .*

- $\partial J_{\rho}(\theta)$ : set of Fréchet sub-differential (Kruger, 2003) of  $J_{\rho}$  at  $\theta$
- Reduces to vanilla policy gradient if  $R = 0$

# Robust Policy Sub-gradient Algorithm

**Input:**  $T, \alpha_t$

**Initialization:**  $\theta_0$

**FOR**  $t = 0, 1, \dots, T - 1$

$$\theta_{t+1} \leftarrow \Pi_{\Theta}(\theta_t - \alpha_t \psi_{\mu}(\theta_t))$$

**Output:**  $\theta$

# Robust Policy Sub-gradient Algorithm

**Input:**  $T, \alpha_t$

**Initialization:**  $\theta_0$

**FOR**  $t = 0, 1, \dots, T - 1$

$$\theta_{t+1} \leftarrow \prod_{\Theta} (\theta_t - \alpha_t \psi_{\mu}(\theta_t))$$

**Output:**  $\theta$

- Vanilla policy gradient is able to find globally optimal policy for non-robust RL, e.g., (Bhandari and Russo, 2021; Agarwal et al., 2021; Cen et al., 2021)

# Robust Policy Sub-gradient Algorithm

**Input:**  $T, \alpha_t$

**Initialization:**  $\theta_0$

**FOR**  $t = 0, 1, \dots, T - 1$

$$\theta_{t+1} \leftarrow \prod_{\Theta} (\theta_t - \alpha_t \psi_{\mu}(\theta_t))$$

**Output:**  $\theta$

- Vanilla policy gradient is able to find globally optimal policy for non-robust RL, e.g., (Bhandari and Russo, 2021; Agarwal et al., 2021; Cen et al., 2021)
- Question: is robust policy sub-gradient able to converge to global optimum of  $J_{\rho}(\theta)$ ?

# Robust Policy Sub-gradient Algorithm

**Input:**  $T, \alpha_t$

**Initialization:**  $\theta_0$

**FOR**  $t = 0, 1, \dots, T - 1$

$$\theta_{t+1} \leftarrow \prod_{\Theta} (\theta_t - \alpha_t \psi_{\mu}(\theta_t))$$

**Output:**  $\theta$

- Vanilla policy gradient is able to find globally optimal policy for non-robust RL, e.g., (Bhandari and Russo, 2021; Agarwal et al., 2021; Cen et al., 2021)
- Question: is robust policy sub-gradient able to converge to global optimum of  $J_{\rho}(\theta)$ ?
- Answer: **Yes!**

PL-condition (Karimi et al., 2016; Bolte et al., 2007):

## Theorem (PL-Condition)

*Under direct policy parameterization,*

$$J_\rho(\theta) - J_\rho^* \leq C_{PL} \max_{\hat{\pi} \in (\Delta(\mathcal{A}))^{|S|}} \langle \pi_\theta - \hat{\pi}, \psi_\rho(\theta) \rangle .$$

## Theorem (Global Optimality under Direct Parameterization)

*If  $\alpha_t > 0$ ,  $\sum_{t=0}^{\infty} \alpha_t = \infty$  and  $\sum_{t=0}^{\infty} \alpha_t^2 < \infty$ , then under direct policy parameterization,  $\theta_T$  converges to a global optimum of  $J_\rho(\theta)$  as  $T \rightarrow \infty$  almost surely.*

## Theorem (Global Optimality under Direct Parameterization)

*If  $\alpha_t > 0$ ,  $\sum_{t=0}^{\infty} \alpha_t = \infty$  and  $\sum_{t=0}^{\infty} \alpha_t^2 < \infty$ , then under direct policy parameterization,  $\theta_T$  converges to a global optimum of  $J_\rho(\theta)$  as  $T \rightarrow \infty$  almost surely.*

- Sub-gradient method converges to stationary points:  $\{\theta : 0 \in \partial J_\rho(\theta)\}$
- Stationary point is globally optimal due to PL-condition



4 Robust Discounted RL

5 Robust Sub-gradient

# Smoothed Robust Policy Gradient

Robust policy sub-gradient method:

- Complexity is generally difficult to establish

# Smoothed Robust Policy Gradient

Robust policy sub-gradient method:

- Complexity is generally difficult to establish

Solution: smoothed robust policy gradient

# Smoothed Robust Policy Gradient

Smoothed robust Bellman operator:

$$\mathbf{T}_{\sigma}^{\pi} V(s) = \mathbb{E}_{A \sim \pi(\cdot|s)} \left[ c(s, A) + \gamma(1 - R) \sum_{s' \in \mathcal{S}} p_{s,s'}^A V(s') + \gamma R \cdot \text{LSE}(\sigma, V) \right],$$

where  $\text{LSE}(\sigma, V) = \frac{\log(\sum_{i=1}^d e^{\sigma V(i)})}{\sigma}$  for  $V \in \mathbb{R}^d$  and some  $\sigma > 0$

# Smoothed Robust Policy Gradient

Smoothed robust Bellman operator:

$$\mathbf{T}_{\sigma}^{\pi} V(s) = \mathbb{E}_{A \sim \pi(\cdot|s)} \left[ c(s, A) + \gamma(1 - R) \sum_{s' \in \mathcal{S}} p_{s,s'}^A V(s') + \gamma R \cdot \text{LSE}(\sigma, V) \right],$$

where  $\text{LSE}(\sigma, V) = \frac{\log(\sum_{i=1}^d e^{\sigma V(i)})}{\sigma}$  for  $V \in \mathbb{R}^d$  and some  $\sigma > 0$

- $\text{LSE}(\sigma, V)$  converges to  $\max_s V(s)$  as  $\sigma \rightarrow \infty$

# Smoothed Robust Policy Gradient

Smoothed robust Bellman operator:

$$\mathbf{T}_\sigma^\pi V(s) = \mathbb{E}_{A \sim \pi(\cdot|s)} \left[ c(s, A) + \gamma(1 - R) \sum_{s' \in \mathcal{S}} p_{s,s'}^A V(s') + \gamma R \cdot \text{LSE}(\sigma, V) \right],$$

where  $\text{LSE}(\sigma, V) = \frac{\log(\sum_{i=1}^d e^{\sigma V(i)})}{\sigma}$  for  $V \in \mathbb{R}^d$  and some  $\sigma > 0$

- $\text{LSE}(\sigma, V)$  converges to  $\max_s V(s)$  as  $\sigma \rightarrow \infty$
- $T_\sigma^\pi$  is a contraction,  $V_\sigma^\pi$  is the fixed point of  $T_\sigma^\pi$   
softmax will not induce contraction (Asadi and Littman, 2017)

# Smoothed Robust Policy Gradient

Smoothed robust Bellman operator:

$$\mathbf{T}_\sigma^\pi V(s) = \mathbb{E}_{A \sim \pi(\cdot|s)} \left[ c(s, A) + \gamma(1 - R) \sum_{s' \in \mathcal{S}} p_{s,s'}^A V(s') + \gamma R \cdot \text{LSE}(\sigma, V) \right],$$

where  $\text{LSE}(\sigma, V) = \frac{\log(\sum_{i=1}^d e^{\sigma V(i)})}{\sigma}$  for  $V \in \mathbb{R}^d$  and some  $\sigma > 0$

- $\text{LSE}(\sigma, V)$  converges to  $\max_s V(s)$  as  $\sigma \rightarrow \infty$
- $T_\sigma^\pi$  is a contraction,  $V_\sigma^\pi$  is the fixed point of  $T_\sigma^\pi$   
softmax will not induce contraction (Asadi and Littman, 2017)
- $V_\sigma^\pi$  is differentiable in  $\theta$  and converges to  $V^\pi$  as  $\sigma \rightarrow \infty$

# Smoothed Robust Policy Gradient

- $J_{\rho}^{\sigma}(\theta) = \sum_{s \in \mathcal{S}} \rho(s) V_{\sigma}^{\pi_{\theta}}(s)$ : smoothed robust objective



# Smoothed Robust Policy Gradient

- $J_\rho^\sigma(\theta) = \sum_{s \in \mathcal{S}} \rho(s) V_\sigma^{\pi_\theta}(s)$ : smoothed robust objective
- Gradient of  $J_\rho^\sigma(\theta)$ :

$$\nabla J_\rho^\sigma(\theta) = B(\rho, \theta) + \frac{\gamma R \sum_{s \in \mathcal{S}} e^{\sigma V_\sigma^{\pi_\theta}(s)} B(s, \theta)}{(1 - \gamma) \sum_{s \in \mathcal{S}} e^{\sigma V_\sigma^{\pi_\theta}(s)}},$$

where  $B(s, \theta) \triangleq \frac{1}{1 - \gamma + \gamma R} \sum_{s' \in \mathcal{S}} d_s^\pi(s') \sum_{a \in \mathcal{A}} \nabla \pi_\theta(a|s') Q_\sigma^{\pi_\theta}(s', a)$ , and  $B(\rho, \theta) \triangleq \mathbb{E}_{S \sim \rho}[B(S, \theta)]$ .

- Smoothed robust policy gradient:  $\theta_{t+1} \leftarrow \Pi_\Theta(\theta_t - \alpha_t \nabla J_\rho^\sigma(\theta))$

# Smoothed Robust Policy Gradient

- $J_\rho^\sigma(\theta) = \sum_{s \in \mathcal{S}} \rho(s) V_\sigma^{\pi_\theta}(s)$ : smoothed robust objective
- Gradient of  $J_\rho^\sigma(\theta)$ :

$$\nabla J_\rho^\sigma(\theta) = B(\rho, \theta) + \frac{\gamma R \sum_{s \in \mathcal{S}} e^{\sigma V_\sigma^{\pi_\theta}(s)} B(s, \theta)}{(1 - \gamma) \sum_{s \in \mathcal{S}} e^{\sigma V_\sigma^{\pi_\theta}(s)}},$$

where  $B(s, \theta) \triangleq \frac{1}{1 - \gamma + \gamma R} \sum_{s' \in \mathcal{S}} d_s^\pi(s') \sum_{a \in \mathcal{A}} \nabla \pi_\theta(a|s') Q_\sigma^{\pi_\theta}(s', a)$ , and  $B(\rho, \theta) \triangleq \mathbb{E}_{S \sim \rho}[B(S, \theta)]$ .

- Smoothed robust policy gradient:  $\theta_{t+1} \leftarrow \Pi_\Theta(\theta_t - \alpha_t \nabla J_\rho^\sigma(\theta))$

Even though gradient is for  $J_\rho^\sigma$ , the algorithm can still find a global optimum of  $J_\rho$  by choosing a large  $\sigma$

Consider direct policy parameterization

## Theorem

For any  $\epsilon > 0$ , set  $\sigma = \mathcal{O}(\epsilon^{-1})$  and  $T = \mathcal{O}(\epsilon^{-3})$ , then

$$\min_{t \leq T-1} J(\theta_t) - J^* \leq 3\epsilon.$$

Consider direct policy parameterization

## Theorem

For any  $\epsilon > 0$ , set  $\sigma = \mathcal{O}(\epsilon^{-1})$  and  $T = \mathcal{O}(\epsilon^{-3})$ , then

$$\min_{t \leq T-1} J(\theta_t) - J^* \leq 3\epsilon.$$

- If  $R = 0$ , i.e., no robustness is considered, complexity reduces to  $\mathcal{O}(\epsilon^{-2})$ , which matches with vanilla policy gradient in (Agarwal et al., 2021)

# Model-free Robust Actor-Critic

- Recall robust policy subgradient:

$$\begin{aligned}\psi_\rho(\theta) \triangleq & \frac{\gamma R}{(1-\gamma)(1-\gamma+\gamma R)} \sum_{s \in \mathcal{S}} d_{s\theta}^{\pi_\theta}(s) \sum_{a \in \mathcal{A}} \nabla \pi_\theta(a|s) Q_{\mathcal{P}, \gamma}^{\pi_\theta}(s, a) \\ & + \frac{1}{1-\gamma+\gamma R} \sum_{s \in \mathcal{S}} d_\rho^{\pi_\theta}(s) \sum_{a \in \mathcal{A}} \nabla \pi_\theta(a|s) Q_{\mathcal{P}, \gamma}^{\pi_\theta}(s, a)\end{aligned}$$

- $Q^{\pi_\theta}(s, a)$  measures cost under worst-case transition kernel and  $\pi_\theta$ , however, only samples from simulator are available

# Model-free Robust Actor-Critic

- Recall robust policy subgradient:

$$\begin{aligned}\psi_\rho(\theta) \triangleq & \frac{\gamma R}{(1-\gamma)(1-\gamma+\gamma R)} \sum_{s \in \mathcal{S}} d_{s\theta}^{\pi_\theta}(s) \sum_{a \in \mathcal{A}} \nabla \pi_\theta(a|s) Q_{\mathcal{P},\gamma}^{\pi_\theta}(s, a) \\ & + \frac{1}{1-\gamma+\gamma R} \sum_{s \in \mathcal{S}} d_\rho^{\pi_\theta}(s) \sum_{a \in \mathcal{A}} \nabla \pi_\theta(a|s) Q_{\mathcal{P},\gamma}^{\pi_\theta}(s, a)\end{aligned}$$

- $Q^{\pi_\theta}(s, a)$  measures cost under worst-case transition kernel and  $\pi_\theta$ , however, only samples from simulator are available

Monte Carlo does not work

# Critic: Robust TD

- Parametric robust action value function  $Q_\zeta$ , e.g., linear function approximation or neural network.

**Input:**  $T_c, \pi, \beta_t$

**Initialization:**  $\zeta, s_0$

Choose  $a_0 \sim \pi(\cdot|s_0)$

**FOR**  $t = 0, 1, \dots, T_c - 1$

Observe  $c_t, s_{t+1}$

Choose  $a_{t+1} \sim \pi(\cdot|s_{t+1})$

$V_t^* \leftarrow \max_s \left\{ \sum_{a \in \mathcal{A}} \pi(a|s) Q_\zeta(s, a) \right\}$

$\delta_t \leftarrow Q_\zeta(s_t, a_t) - \underbrace{(c_t + \gamma(1 - R)Q_\zeta(s_{t+1}, a_{t+1}) + \gamma RV_t^*)}_{\text{robust target}} \text{ (robust TD error)}$

$\zeta \leftarrow \zeta - \beta_t \delta_t \nabla_\zeta Q_\zeta(s_t, a_t)$

**Output:**  $\zeta$

# Robust Actor-Critic Algorithm

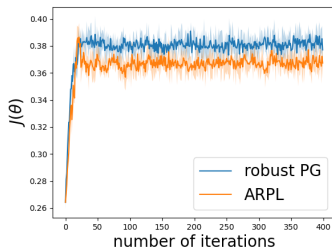
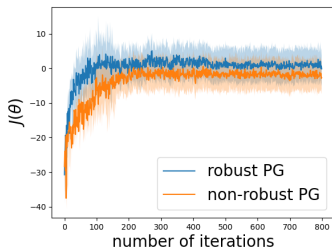
- Using robust TD algorithm to estimate robust Q-function in (smoothed) robust policy gradient
- Under tabular setting, global optimality can be established, overall sample complexity is  $\mathcal{O}(\epsilon^{-7})$

Robust actor-critic algorithm can be applied with arbitrary value function/policy approximation.



# Experiments

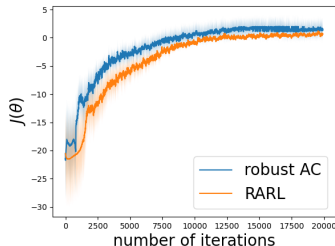
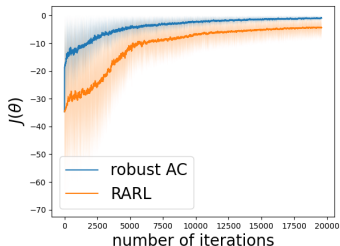
- Robust policy gradient v.s. vanilla policy gradient and ARPL  
Mandlekar et al. (2017)
- ARPL: Adversary randomly perturb observation then run vanilla policy gradient method using these perturbed samples
- Training on an unperturbed MDP, and evaluation on the worst-case transition kernel in  $\mathcal{P}$



- Our robust policy gradient achieves higher reward on the worst-case transition kernel

# Experiments

- Robust actor-critic v.s. RARL (Pinto et al., 2017)
- RARL: Adversary perturbs state transition. Agent and adversary are updated alternatively using gradient descent ascent.
- Training on an unperturbed MDP, and evaluation on the worst-case transition kernel in  $\mathcal{P}$



- Our robust actor critic achieves higher reward on the worst-case transition kernel

# Summary

- Robust policy gradient with provable global optimality
- Model-free robust actor-critic algorithm
- Can be easily scaled to large/continuous problems