

Model-free Approaches to Robust Markov Decision Processes

Wenhao Yang

Stanford Univeristy

2023.10.17

- ① Introduction
- ② Model-free robust MDPs
- ③ Conclusion
- ④ Reference

1 Introduction

2 Model-free robust MDPs

3 Conclusion

4 Reference

What is Robust RL?

- Two environments: A and B.
- Optimal policy in A may be sub-optimal/bad in B.
- Can we design a conservative policy?

What is Robust MDPs?

- Instance parameters: $\langle \mathcal{S}, \mathcal{A}, R, P^*, \gamma \rangle$.

What is Robust MDPs?

- Instance parameters: $\langle \mathcal{S}, \mathcal{A}, R, P^*, \gamma \rangle$.
- Uncertainty set: $\mathcal{P} = \bigotimes_{s,a} \mathcal{P}_{s,a}$,

What is Robust MDPs?

- Instance parameters: $\langle \mathcal{S}, \mathcal{A}, R, P^*, \gamma \rangle$.
- Uncertainty set: $\mathcal{P} = \bigotimes_{s,a} \mathcal{P}_{s,a}$,

$$\mathcal{P}_{s,a} := \left\{ P(\cdot|s, a) \in \Delta(\mathcal{S}) \mid \sum_{s' \in \mathcal{S}} f \left(\frac{P(s'|s, a)}{P^*(s'|s, a)} \right) P^*(s'|s, a) \leq \rho \right\}.$$

What is Robust MDPs?

- Instance parameters: $\langle \mathcal{S}, \mathcal{A}, R, P^*, \gamma \rangle$.
- Uncertainty set: $\mathcal{P} = \bigotimes_{s,a} \mathcal{P}_{s,a}$,

$$\mathcal{P}_{s,a} := \left\{ P(\cdot|s,a) \in \Delta(\mathcal{S}) \mid \sum_{s' \in \mathcal{S}} f \left(\frac{P(s'|s,a)}{P^*(s'|s,a)} \right) P^*(s'|s,a) \leq \rho \right\}.$$

- Robust value function with given policy π :

What is Robust MDPs?

- Instance parameters: $\langle \mathcal{S}, \mathcal{A}, R, P^*, \gamma \rangle$.
- Uncertainty set: $\mathcal{P} = \bigotimes_{s,a} \mathcal{P}_{s,a}$,

$$\mathcal{P}_{s,a} := \left\{ P(\cdot|s,a) \in \Delta(\mathcal{S}) \mid \sum_{s' \in \mathcal{S}} f \left(\frac{P(s'|s,a)}{P^*(s'|s,a)} \right) P^*(s'|s,a) \leq \rho \right\}.$$

- Robust value function with given policy π :

$$V_{\text{rob,c}}^\pi(s) := \inf_{P \in \mathcal{P}} V_P^\pi(s).$$

What is Robust MDPs?

- Instance parameters: $\langle \mathcal{S}, \mathcal{A}, R, P^*, \gamma \rangle$.
- Uncertainty set: $\mathcal{P} = \bigotimes_{s,a} \mathcal{P}_{s,a}$,

$$\mathcal{P}_{s,a} := \left\{ P(\cdot|s,a) \in \Delta(\mathcal{S}) \mid \sum_{s' \in \mathcal{S}} f \left(\frac{P(s'|s,a)}{P^*(s'|s,a)} \right) P^*(s'|s,a) \leq \rho \right\}.$$

- Robust value function with given policy π :

$$V_{\text{rob},c}^{\pi}(s) := \inf_{P \in \mathcal{P}} V_P^{\pi}(s).$$

- Optimal robust value function $V_{\text{rob},c}^*(s) := \max_{\pi} V_{\text{rob},c}^{\pi}(s)$.

Why Robust MDPs?

- Why care robust MDPs?

Why Robust MDPs?

- Why care robust MDPs?
- A large ρ can reduce sample complexity. [YZZ22]

Why Robust MDPs?

- Why care robust MDPs?
- A large ρ can reduce sample complexity. [YZZ22]

Theorem [YZZ22]

There exists a class of MDPs, given that $f(t) = (t - 1)^2$, for every (ε, δ) -correct RL algorithm \mathcal{A} , the total number of samples needs to be at least:

$$\tilde{\Omega} \left(\frac{|\mathcal{S}||\mathcal{A}|}{\varepsilon^2(1-\gamma)^2} \min \left\{ \frac{1}{1-\gamma}, \frac{1}{\rho} \right\} \right). \quad (1)$$

How to solve Robust MDPs?

- $V_{rob,c}^*$ still satisfies Bellman equation:

How to solve Robust MDPs?

- $V_{\text{rob,c}}^*$ still satisfies Bellman equation:

$$V_{\text{rob,c}}^*(s) = \max_a \left(R(s, a) + \gamma \inf_{P_{s,a} \in \mathcal{P}_{s,a}} \sum_{s'} P(s'|s, a) V_{\text{rob,c}}^*(s') \right)$$
$$:= \mathcal{T}_{\text{rob,c}} V_{\text{rob,c}}^*(s).$$

How to solve Robust MDPs?

- $V_{\text{rob},c}^*$ still satisfies Bellman equation:

$$V_{\text{rob},c}^*(s) = \max_a \left(R(s, a) + \gamma \inf_{P_{s,a} \in \mathcal{P}_{s,a}} \sum_{s'} P(s'|s, a) V_{\text{rob},c}^*(s') \right) \\ := \mathcal{T}_{\text{rob},c} V_{\text{rob},c}^*(s).$$

- If we know:
 - A good estimation of P^* ,
 - Solution of inner optimization problem $\inf_{P_{s,a} \in \mathcal{P}_{s,a}} \sum_{s'} P(s'|s, a) V(s')$ for any given V ,

How to solve Robust MDPs?

- $V_{\text{rob},c}^*$ still satisfies Bellman equation:

$$V_{\text{rob},c}^*(s) = \max_a \left(R(s, a) + \gamma \inf_{P_{s,a} \in \mathcal{P}_{s,a}} \sum_{s'} P(s'|s, a) V_{\text{rob},c}^*(s') \right)$$
$$:= \mathcal{T}_{\text{rob},c} V_{\text{rob},c}^*(s).$$

- If we know:
 - A good estimation of P^* ,
 - Solution of inner optimization problem $\inf_{P_{s,a} \in \mathcal{P}_{s,a}} \sum_{s'} P(s'|s, a) V(s')$ for any given V ,near-optimal robust value function can be obtained with efficient sample complexity. [YZZ22]

Problems on Current Setting.

- Memory space $\mathcal{O}(|\mathcal{S}|^2|\mathcal{A}|)$. (Storing \hat{P})

Problems on Current Setting.

- Memory space $\mathcal{O}(|\mathcal{S}|^2|\mathcal{A}|)$. (Storing \hat{P})
- Computation complexity of inner optimization problem enlarges with instance size.

Problems on Current Setting.

- Memory space $\mathcal{O}(|\mathcal{S}|^2|\mathcal{A}|)$. (Storing \hat{P})
- Computation complexity of inner optimization problem enlarges with instance size.
- Question: can we design a model-free algorithm with efficient sample complexity (including the complexity of inner optimization problem)?

- ① Introduction
- ② Model-free robust MDPs**
- ③ Conclusion
- ④ Reference

Q-learning

- Q-learning updating rule:

Q-learning

- Q-learning updating rule:

$$Q_{t+1}(s, a) = (1 - \alpha_t)Q_t(s, a) + \alpha_t \hat{\mathcal{T}}Q_t(s, a), \quad (2)$$

Q-learning

- Q-learning updating rule:

$$Q_{t+1}(s, a) = (1 - \alpha_t)Q_t(s, a) + \alpha_t \widehat{\mathcal{T}}Q_t(s, a), \quad (2)$$

where $\widehat{\mathcal{T}}Q_t(s, a) = r_t(s, a) + \gamma \max_a Q_t(s'_t, a)$

Q-learning

- Q-learning updating rule:

$$Q_{t+1}(s, a) = (1 - \alpha_t)Q_t(s, a) + \alpha_t \widehat{\mathcal{T}}Q_t(s, a), \quad (2)$$

where $\widehat{\mathcal{T}}Q_t(s, a) = r_t(s, a) + \gamma \max_a Q_t(s'_t, a)$ and satisfies $\mathbb{E}[\widehat{\mathcal{T}}Q_t | \mathcal{F}_t] = \mathcal{T}Q_t = R(s, a) + \gamma \mathbb{E}_P \max_{a'} Q_t(s', a)$.

Q-learning

- Q-learning updating rule:

$$Q_{t+1}(s, a) = (1 - \alpha_t)Q_t(s, a) + \alpha_t \widehat{\mathcal{T}}Q_t(s, a), \quad (2)$$

where $\widehat{\mathcal{T}}Q_t(s, a) = r_t(s, a) + \gamma \max_a Q_t(s'_t, a)$ and satisfies $\mathbb{E}[\widehat{\mathcal{T}}Q_t | \mathcal{F}_t] = \mathcal{T}Q_t = R(s, a) + \gamma \mathbb{E}_P \max_{a'} Q_t(s', a)$.

- Subtracting Q^* each side:

Q-learning

- Q-learning updating rule:

$$Q_{t+1}(s, a) = (1 - \alpha_t)Q_t(s, a) + \alpha_t \widehat{\mathcal{T}}Q_t(s, a), \quad (2)$$

where $\widehat{\mathcal{T}}Q_t(s, a) = r_t(s, a) + \gamma \max_a Q_t(s'_t, a)$ and satisfies $\mathbb{E}[\widehat{\mathcal{T}}Q_t | \mathcal{F}_t] = \mathcal{T}Q_t = R(s, a) + \gamma \mathbb{E}_P \max_{a'} Q_t(s', a)$.

- Subtracting Q^* each side:

$$\begin{aligned} Q_{t+1} - Q^* &= (1 - \alpha_t)(Q_t - Q^*) \\ &\quad + \alpha_t(\widehat{\mathcal{T}}Q_t - \mathcal{T}Q_t) + \alpha_t(\mathcal{T}Q_t - \mathcal{T}Q^*) \end{aligned} \quad (3)$$

Theorem [Wai19, LYZJ21]

For a sequence $X_{t+1} = (1 - \alpha_t)X_t + \alpha_t Y_t + \alpha_t \delta_t$, where $\{Y_t\}_{t \geq 0}$ is a martingale difference, $(1 - \alpha_t)\alpha_{t-1} \leq \alpha_t$, and $\sum_{t=0}^{T-1} \delta_t = o(1/\alpha_T)$, then X_T satisfies:

$$\mathbb{E}|X_T| \leq \tilde{O} \left(\sqrt{\alpha_T} + \alpha_T \sum_{t=0}^{T-1} \delta_t \right) \quad (4)$$

Theorem [Wai19, LYZJ21]

For a sequence $X_{t+1} = (1 - \alpha_t)X_t + \alpha_t Y_t + \alpha_t \delta_t$, where $\{Y_t\}_{t \geq 0}$ is a martingale difference, $(1 - \alpha_t)\alpha_{t-1} \leq \alpha_t$, and $\sum_{t=0}^{T-1} \delta_t = o(1/\alpha_T)$, then X_T satisfies:

$$\mathbb{E}|X_T| \leq \tilde{\mathcal{O}} \left(\sqrt{\alpha_T} + \alpha_T \sum_{t=0}^{T-1} \delta_t \right) \quad (4)$$

- Convergence rate is guaranteed [Wai19]:

$$\mathbb{E} \|Q_T - Q^*\|_\infty = \tilde{\mathcal{O}}(\sqrt{\alpha_T}) \quad (5)$$

- For robust MDPs, noting the robust Bellman operator:

$$\mathcal{T}_{\text{rob},c}Q = R(s, a) + \gamma \inf_{D_f(P \| P_{s,a}^*) \leq \rho} \mathbb{E}_{s' \sim P} \max_a Q(s', a) \quad (6)$$

- For robust MDPs, noting the robust Bellman operator:

$$\mathcal{T}_{\text{rob},c}Q = R(s, a) + \gamma \inf_{D_f(P \| P_{s,a}^*) \leq \rho} \mathbb{E}_{s' \sim P} \max_a Q(s', a) \quad (6)$$

- Non-linear functional of expectation.

Key observation

- Simplify the notation of inner optimization problem:

Key observation

- Simplify the notation of inner optimization problem:

$$\inf_{D_f(P\|P^*) \leq \rho} \sum_i P_i V_i.$$

Key observation

- Simplify the notation of inner optimization problem:

$$\inf_{D_f(P\|P^*) \leq \rho} \sum_i P_i V_i.$$

- Dual form:

$$\sup_{\lambda \geq 0, \eta \in \mathbb{R}} -\lambda \sum_i P_i^* f^* \left(\frac{\eta - V_i}{\lambda} \right) - \lambda \rho + \eta,$$

where $f^*(t) = \sup_{s \geq 0} (st - f(s))$.

Key observation

- Simplify the notation of inner optimization problem:

$$\inf_{D_f(P\|P^*) \leq \rho} \sum_i P_i V_i.$$

- Dual form:

$$\sup_{\lambda \geq 0, \eta \in \mathbb{R}} -\lambda \sum_i P_i^* f^* \left(\frac{\eta - V_i}{\lambda} \right) - \lambda \rho + \eta,$$

where $f^*(t) = \sup_{s \geq 0} (st - f(s))$.

- How to construct a good estimator for $\mathcal{T}_r V$ with a given V and $\mathcal{O}(1)$ samples?
(specifically, unbiased)

Key observation

- Problem is equivalent to construct a good estimator for:

Key observation

- Problem is equivalent to construct a good estimator for:

$$\sup_{\theta} \mathbb{E}_P f(X; \theta). \quad (7)$$

Key observation

- Problem is equivalent to construct a good estimator for:

$$\sup_{\theta} \mathbb{E}_P f(X; \theta). \quad (7)$$

- But $\mathbb{E} \sup_{\theta} \frac{1}{n} \sum_{i=1}^m f(X_i; \theta) \neq \sup_{\theta} \mathbb{E}_P f(X; \theta)$.

Key observation

- Problem is equivalent to construct a good estimator for:

$$\sup_{\theta} \mathbb{E}_P f(X; \theta). \quad (7)$$

- But $\mathbb{E} \sup_{\theta} \frac{1}{n} \sum_{i=1}^m f(X_i; \theta) \neq \sup_{\theta} \mathbb{E}_P f(X; \theta)$.
- Can we construct a random variable Z_n based on $\{X_i\}$ s.t. $\mathbb{E} Z_n = \sup_{\theta} \mathbb{E}_P f(X; \theta)$?

Multilevel Monte-Carlo method

- Multilevel Monte-Carlo method [BGP19]: Given 2^{N+1} i.i.d. samples $\{X_i\}_{i=1}^{2^{N+1}}$ with $N \sim \text{Geo}(g)$,

Multilevel Monte-Carlo method

- Multilevel Monte-Carlo method [BGP19]: Given 2^{N+1} i.i.d. samples $\{X_i\}_{i=1}^{2^{N+1}}$ with $N \sim \text{Geo}(g)$,

$$\begin{aligned} \Delta_N &= \sup_{\theta} \frac{1}{2^{N+1}} \sum_{i=1}^{2^{N+1}} f(X_i; \theta) \\ &\quad - \frac{1}{2} \sup_{\theta} \frac{1}{2^N} \sum_{i=1}^{2^N} f(X_{2i}; \theta) - \frac{1}{2} \sup_{\theta} \frac{1}{2^N} \sum_{i=1}^{2^N} f(X_{2i-1}; \theta). \end{aligned} \quad (8)$$

Multilevel Monte-Carlo method

- Multilevel Monte-Carlo method [BGP19]: Given 2^{N+1} i.i.d. samples $\{X_i\}_{i=1}^{2^{N+1}}$ with $N \sim \text{Geo}(g)$,

$$\begin{aligned} \Delta_N = & \sup_{\theta} \frac{1}{2^{N+1}} \sum_{i=1}^{2^{N+1}} f(X_i; \theta) \\ & - \frac{1}{2} \sup_{\theta} \frac{1}{2^N} \sum_{i=1}^{2^N} f(X_{2i}; \theta) - \frac{1}{2} \sup_{\theta} \frac{1}{2^N} \sum_{i=1}^{2^N} f(X_{2i-1}; \theta). \end{aligned} \quad (8)$$

- $\mathbb{E}[\sup_{\theta} f(X_1; \theta) + \Delta_N/p_N] = \sup_{\theta} \mathbb{E}_P f(X; \theta)$ by noticing:

Multilevel Monte-Carlo method

- Multilevel Monte-Carlo method [BGP19]: Given 2^{N+1} i.i.d. samples $\{X_i\}_{i=1}^{2^{N+1}}$ with $N \sim \text{Geo}(g)$,

$$\begin{aligned} \Delta_N &= \sup_{\theta} \frac{1}{2^{N+1}} \sum_{i=1}^{2^{N+1}} f(X_i; \theta) \\ &\quad - \frac{1}{2} \sup_{\theta} \frac{1}{2^N} \sum_{i=1}^{2^N} f(X_{2i}; \theta) - \frac{1}{2} \sup_{\theta} \frac{1}{2^N} \sum_{i=1}^{2^N} f(X_{2i-1}; \theta). \end{aligned} \quad (8)$$

- $\mathbb{E}[\sup_{\theta} f(X_1; \theta) + \Delta_N/p_N] = \sup_{\theta} \mathbb{E}_P f(X; \theta)$ by noticing:

$$\mathbb{E}[\Delta_N/p_N] = \sum_{n=0}^{+\infty} \mathbb{E}_{P_{2^{n+1}}} [\sup_{\theta} f(X; \theta)] - \mathbb{E}_{P_{2^n}} [\sup_{\theta} f(X; \theta)] \quad (9)$$

Distributionally Robust Q-learning

- Updating rule of DR Q-learning [LBB⁺22]:

Distributionally Robust Q-learning

- Updating rule of DR Q-learning [LBB⁺22]:
 - For each (s, a) , sample $N \sim \text{Geo}(g)$, and 2^{N+1} samples of $s' \sim P(\cdot|s, a)$.

Distributionally Robust Q-learning

- Updating rule of DR Q-learning [LBB⁺22]:
 - For each (s, a) , sample $N \sim \text{Geo}(g)$, and 2^{N+1} samples of $s' \sim P(\cdot|s, a)$.
 - $Q_{t+1}(s, a) = (1 - \alpha_t)Q_t(s, a) + \alpha_t \widehat{\mathcal{T}}_{\text{rob},c} Q_t(s, a)$, where $\widehat{\mathcal{T}}_{\text{rob},c} Q_t$ is the Multilevel MC estimator.

Distributionally Robust Q-learning

- Updating rule of DR Q-learning [LBB⁺22]:
 - For each (s, a) , sample $N \sim \text{Geo}(g)$, and 2^{N+1} samples of $s' \sim P(\cdot|s, a)$.
 - $Q_{t+1}(s, a) = (1 - \alpha_t)Q_t(s, a) + \alpha_t \widehat{\mathcal{T}}_{\text{rob},c} Q_t(s, a)$, where $\widehat{\mathcal{T}}_{\text{rob},c} Q_t$ is the Multilevel MC estimator.
- Choice of parameter $g \in (1/2, 3/4)$, leading $\mathbb{E}2^{N+1} = 2g/(2g - 1)$.

Distributionally Robust Q-learning

- Updating rule of DR Q-learning [LBB⁺22]:
 - For each (s, a) , sample $N \sim \text{Geo}(g)$, and 2^{N+1} samples of $s' \sim P(\cdot|s, a)$.
 - $Q_{t+1}(s, a) = (1 - \alpha_t)Q_t(s, a) + \alpha_t \widehat{\mathcal{T}}_{\text{rob},c} Q_t(s, a)$, where $\widehat{\mathcal{T}}_{\text{rob},c} Q_t$ is the Multilevel MC estimator.
- Choice of parameter $g \in (1/2, 3/4)$, leading $\mathbb{E}2^{N+1} = 2g/(2g - 1)$.
- Sample complexity [WSBZ23]: $\tilde{O}\left(\frac{|S||A|}{\varepsilon^2(1-\gamma)^5\rho^4}\right)$.

Distributionally Robust Q-learning

- Updating rule of DR Q-learning [LBB⁺22]:
 - For each (s, a) , sample $N \sim \text{Geo}(g)$, and 2^{N+1} samples of $s' \sim P(\cdot|s, a)$.
 - $Q_{t+1}(s, a) = (1 - \alpha_t)Q_t(s, a) + \alpha_t \widehat{\mathcal{T}}_{\text{rob},c} Q_t(s, a)$, where $\widehat{\mathcal{T}}_{\text{rob},c} Q_t$ is the Multilevel MC estimator.
- Choice of parameter $g \in (1/2, 3/4)$, leading $\mathbb{E}2^{N+1} = 2g/(2g - 1)$.
- Sample complexity [WSBZ23]: $\tilde{O}\left(\frac{|S||A|}{\varepsilon^2(1-\gamma)^5\rho^4}\right)$.
- Limitation: unknown to computation complexity.

Alternative form

- Recall the dual form:

Alternative form

- Recall the dual form:

$$\sup_{\lambda \geq 0, \eta \in \mathbb{R}} -\lambda \sum_i P_i^* f^* \left(\frac{\eta - V_i}{\lambda} \right) - \lambda \rho + \eta.$$

Alternative form

- Recall the dual form:

$$\sup_{\lambda \geq 0, \eta \in \mathbb{R}} -\lambda \sum_i P_i^* f^* \left(\frac{\eta - V_i}{\lambda} \right) - \lambda \rho + \eta.$$

- Can we apply stochastic gradient method?

Alternative form

- Recall the dual form:

$$\sup_{\lambda \geq 0, \eta \in \mathbb{R}} -\lambda \sum_i P_i^* f^* \left(\frac{\eta - V_i}{\lambda} \right) - \lambda \rho + \eta.$$

- Can we apply stochastic gradient method?
- Sadly, stochastic gradient method will fail here. [ND16]

Alternative form

- Recall the dual form:

$$\sup_{\lambda \geq 0, \eta \in \mathbb{R}} -\lambda \sum_i P_i^* f^* \left(\frac{\eta - V_i}{\lambda} \right) - \lambda \rho + \eta.$$

- Can we apply stochastic gradient method?
- Sadly, stochastic gradient method will fail here. [ND16]
- Reason: no bounded gradient because of λ .

Alternative form

- Recall the dual form:

$$\sup_{\lambda \geq 0, \eta \in \mathbb{R}} -\lambda \sum_i P_i^* f^* \left(\frac{\eta - V_i}{\lambda} \right) - \lambda \rho + \eta.$$

- Can we apply stochastic gradient method?
- Sadly, stochastic gradient method will fail here. [ND16]
- Reason: no bounded gradient because of λ .
- Set $f(t) = (t - 1)^2$, dual form:

$$\sup_{\lambda \geq 0, \eta \in \mathbb{R}} -\frac{\sum_i P_i^* (\eta - V_i)_+^2}{4\lambda} - \lambda \rho + \eta - \lambda. \quad (10)$$

Alternative form

- Recall the dual form:

$$\sup_{\lambda \geq 0, \eta \in \mathbb{R}} -\lambda \sum_i P_i^* f^* \left(\frac{\eta - V_i}{\lambda} \right) - \lambda \rho + \eta.$$

- Can we apply stochastic gradient method?
- Sadly, stochastic gradient method will fail here. [ND16]
- Reason: no bounded gradient because of λ .
- Set $f(t) = (t - 1)^2$, dual form:

$$\sup_{\lambda \geq 0, \eta \in \mathbb{R}} -\frac{\sum_i P_i^* (\eta - V_i)_+^2}{4\lambda} - \lambda \rho + \eta - \lambda. \quad (10)$$

- Jointly optimizing over λ, η fails.

Alternative form

- What if fix λ ?

Alternative form

- What if fix λ ?
- Let's try penalty form:

$$\inf_P \sum_i P_i V_i + \lambda D_f(P \| P^*)$$

Alternative form

- What if fix λ ?
- Let's try penalty form:

$$\inf_P \sum_i P_i V_i + \lambda D_f(P \| P^*)$$

- Its dual form satisfies:

$$\sup_{\eta \in \mathbb{R}} -\lambda \sum_i P_i^* f^* \left(\frac{\eta - V_i}{\lambda} \right) + \eta. \quad (11)$$

Alternative form

- What if fix λ ?
- Let's try penalty form:

$$\inf_P \sum_i P_i V_i + \lambda D_f(P \| P^*)$$

- Its dual form satisfies:

$$\sup_{\eta \in \mathbb{R}} -\lambda \sum_i P_i^* f^* \left(\frac{\eta - V_i}{\lambda} \right) + \eta. \quad (11)$$

- High-level idea: small ρ (constraint form) corresponds to large λ (penalty form).

Alternative form

- What if fix λ ?
- Let's try penalty form:

$$\inf_P \sum_i P_i V_i + \lambda D_f(P \| P^*)$$

- Its dual form satisfies:

$$\sup_{\eta \in \mathbb{R}} -\lambda \sum_i P_i^* f^* \left(\frac{\eta - V_i}{\lambda} \right) + \eta. \quad (11)$$

- High-level idea: small ρ (constraint form) corresponds to large λ (penalty form).
- Stochastic gradient method works here!

Penalized Robust MDPs

- How to define penalty in robust MDPs?

Penalized Robust MDPs

- How to define penalty in robust MDPs?
- Intuitively, consider a new robust Bellman operator:

$$\mathcal{T}_{\text{rob,p}}V(s) = \max_a \left(R(s, a) + \gamma \inf_P \sum_{s'} P(s')V(s') + \lambda D_f(P \| P_{s,a}^*) \right) \quad (12)$$

Penalized Robust MDPs

- How to define penalty in robust MDPs?
- Intuitively, consider a new robust Bellman operator:

$$\mathcal{T}_{\text{rob,p}}V(s) = \max_a \left(R(s, a) + \gamma \inf_P \sum_{s'} P(s')V(s') + \lambda D_f(P \| P_{s,a}^*) \right) \quad (12)$$

- Define the value function:

$$V_{\text{rob,p}}^\pi(s) := \inf_P \mathbb{E}_{P,\pi} \left[\sum_{t \geq 0} \gamma^t (R(s_t, a_t) + \lambda \gamma D_f(P_{s_t, a_t} \| P_{s_t, a_t}^*)) \mid s_0 = s \right] \quad (13)$$

Penalized Robust MDPs

- Is $V_{\text{rob},p}^\pi$ well defined?

Penalized Robust MDPs

- Is $V_{\text{rob},p}^\pi$ well defined?

Proposition ([YWK⁺23])

$\mathcal{T}_{\text{rob},p}$ is a γ -contraction operator on value space with a fixed point $V_{\text{rob},p}^*$, which satisfies

$$V_{\text{rob},p}^* = \max_{\pi} V_{\text{rob},p}^{\pi}.$$

- Is robustness preserved?

- Is robustness preserved?

Theorem (Statistical Equivalence [YWK⁺23])

With a generative model and $f(t) = (t - 1)^2$, with probability $1 - \delta$,

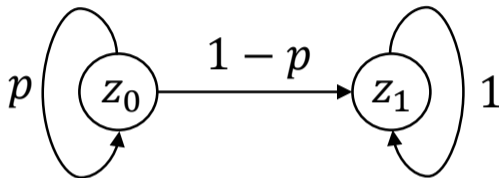
$$\|\widehat{V}_{\text{rob,p}}^* - V_{\text{rob,p}}^*\|_{\infty} \leq \varepsilon, \quad (14)$$

by taking $n = \tilde{\mathcal{O}}\left(\frac{|\mathcal{S}||\mathcal{A}|}{\varepsilon^2(1-\gamma)^2} \max\{\lambda^{-2}(1-\gamma)^{-2}, \lambda^2\}\right)$. Also, there exists a class of robust MDPs, for every (ε, δ) -correct robust RL algorithm, the total number of samples needed is at least:

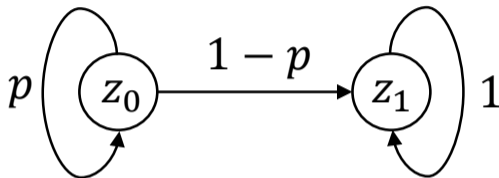
$$n = \begin{cases} \tilde{\Omega}\left(\frac{|\mathcal{S}||\mathcal{A}|\lambda^2}{\varepsilon^2(1-\gamma)^3}\right) & , \text{ when } \lambda = \mathcal{O}(1-\gamma) \\ \tilde{\Omega}\left(\frac{|\mathcal{S}||\mathcal{A}|}{\varepsilon^2(1-\gamma)^3}\right) & , \text{ when } \lambda = \Omega(1-\gamma) \end{cases} \quad (15)$$



- Consider a 2-state MDP:



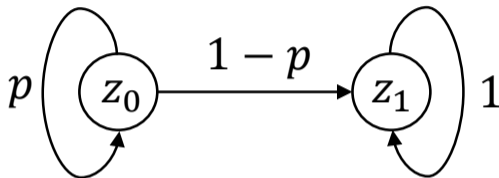
- Consider a 2-state MDP:



- Value function satisfies:

$$V(z_0) = 1 + \gamma \inf_{0 \leq q \leq 1} qV(z_0) + \lambda D_f(q||p) \quad (16)$$

- Consider a 2-state MDP:

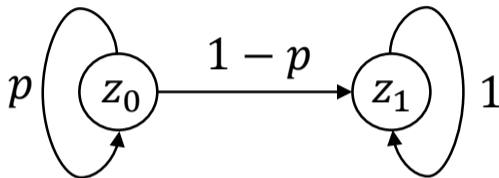


- Value function satisfies:

$$V(z_0) = 1 + \gamma \inf_{0 \leq q \leq 1} qV(z_0) + \lambda D_f(q||p) \quad (16)$$

- For $f(t) = (t - 1)^2$, it is a quadratic equation.

- Consider a 2-state MDP:



- Value function satisfies:

$$V(z_0) = 1 + \gamma \inf_{0 \leq q \leq 1} qV(z_0) + \lambda D_f(q||p) \quad (16)$$

- For $f(t) = (t - 1)^2$, it is a quadratic equation.
- Compare $V(z_0)$ under p and $p + \delta$, then apply information theory.

- We can design the robust Q-learning algorithm with a generative model:

- We can design the robust Q-learning algorithm with a generative model:
At timestep t , for each (s, a) , we do

- We can design the robust Q-learning algorithm with a generative model:
At timestep t , for each (s, a) , we do
 - Stochastic gradient method for $\sup_{\eta \in \mathbb{R}} -\lambda \sum_{s'} P^*(s'|s, a) f^* \left(\frac{\eta - V_t(s')}{\lambda} \right) + \eta$ with sufficient steps T' and obtain $\eta_{T'}(s, a)$.

- We can design the robust Q-learning algorithm with a generative model:
At timestep t , for each (s, a) , we do
 - Stochastic gradient method for $\sup_{\eta \in \mathbb{R}} -\lambda \sum_{s'} P^*(s'|s, a) f^* \left(\frac{\eta - V_t(s')}{\lambda} \right) + \eta$ with sufficient steps T' and obtain $\eta_{T'}(s, a)$.
 - Run one step Q-learning:

$$Q_{t+1}(s, a) = (1 - \beta_t)Q_t(s, a) + \beta_t \widehat{\mathcal{T}}_{\text{rob}, p} V_t,$$

where $\widehat{\mathcal{T}}_{\text{rob}, p} V_t = r_t + \gamma \cdot \left(-\lambda f^* \left(\frac{\eta_{T'}(s, a) - V_t(s')}{\lambda} \right) + \eta_{T'}(s, a) \right)$.

Robust Q-learning

- Why sufficient steps T' for inner optimization problem?

Robust Q-learning

- Why sufficient steps T' for inner optimization problem?
- Given V_t , $|\eta_{T'} - \eta^*|$ is small with high probability.

Robust Q-learning

- Why sufficient steps T' for inner optimization problem?
- Given V_t , $|\eta_{T'} - \eta^*|$ is small with high probability.
- Convergence of Q-learning requires a nearly-unbiased estimator of $\sup_{\eta \in \mathbb{R}} -\lambda \sum_{s'} P^*(s'|s, a) f^* \left(\frac{\eta - V_t(s')}{\lambda} \right) + \eta$.

Robust Q-learning

- Why sufficient steps T' for inner optimization problem?
- Given V_t , $|\eta_{T'} - \eta^*|$ is small with high probability.
- Convergence of Q-learning requires a nearly-unbiased estimator of $\sup_{\eta \in \mathbb{R}} -\lambda \sum_{s'} P^*(s'|s, a) f^* \left(\frac{\eta - V_t(s')}{\lambda} \right) + \eta$.
- With a near optimal $\eta_{T'}$, $-\lambda f^* \left(\frac{\eta_{T'}(s, a) - V_t(s')}{\lambda} \right) + \eta_{T'}(s, a)$ is nearly-unbiased.

Robust Q-learning

- Why sufficient steps T' for inner optimization problem?
- Given V_t , $|\eta_{T'} - \eta^*|$ is small with high probability.
- Convergence of Q-learning requires a nearly-unbiased estimator of $\sup_{\eta \in \mathbb{R}} -\lambda \sum_{s'} P^*(s'|s, a) f^* \left(\frac{\eta - V_t(s')}{\lambda} \right) + \eta$.
- With a near optimal $\eta_{T'}$, $-\lambda f^* \left(\frac{\eta_{T'}(s, a) - V_t(s')}{\lambda} \right) + \eta_{T'}(s, a)$ is nearly-unbiased.
- Error decomposition:

$$\begin{aligned}
 J(s'_t; \eta_{T'}, V_t) - \sup_{\eta} \mathbb{E} J(s'; \eta, V^*) &= J(s'_t; \eta_{T'}, V_t) - \mathbb{E} J(s'; \eta_{T'}, V_t) \\
 &\quad + \mathbb{E} J(s'; \eta_{T'}, V_t) - \sup_{\eta} \mathbb{E} J(s'; \eta, V_t) \\
 &\quad + \sup_{\eta} \mathbb{E} J(s'; \eta, V_t) - \sup_{\eta} \mathbb{E} J(s'; \eta, V^*) \tag{17}
 \end{aligned}$$

Robust Q-learning

Theorem ([YWK⁺23])

Setting $f(t) = (t - 1)^2$, $\alpha_{t'} = \frac{\lambda C}{\sqrt{t'}}$, and $\beta_t = \frac{1}{1+(1-\gamma)(t+1)}$. To obtain an ε -optimal Q-value function, the total number of sample complexity is:

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{\varepsilon^2(1-\gamma)^5}\right) \cdot \tilde{O}\left(\frac{\max\{\lambda^2, \lambda^{-2}(1-\gamma)^{-4}\}}{\varepsilon^2(1-\gamma)^2}\right) \quad (18)$$

Limitations

- [LBB⁺22, WSBZ23, YWK⁺23] require generative model.

Limitations

- [LBB⁺22, WSBZ23, YWK⁺23] require generative model.
- Practical scenario: sampling from one trajectory.

Limitations

- [LBB⁺22, WSBZ23, YWK⁺23] require generative model.
- Practical scenario: sampling from one trajectory.
- Observed samples $(s_0, a_0, r_0, s_1, a_1, \dots)$.

One trajectory

- [LMB⁺23] consider concrete cases: χ^2 and KL.

One trajectory

- [LMB⁺23] consider concrete cases: χ^2 and KL.
- For χ^2 case, optimize over λ :

One trajectory

- [LMB⁺23] consider concrete cases: χ^2 and KL.
- For χ^2 case, optimize over λ :

$$\mathcal{T}_{\text{rob,c}}V = R(s, a) + \gamma \sup_{\eta} \left(\eta - \sqrt{1 + \rho} \sqrt{\mathbb{E}_{P_{s,a}^*} (\eta - V(s'))_+^2} \right) \quad (19)$$

One trajectory

- [LMB⁺23] consider concrete cases: χ^2 and KL.
- For χ^2 case, optimize over λ :

$$\mathcal{T}_{\text{rob,c}}V = R(s, a) + \gamma \sup_{\eta} \left(\eta - \sqrt{1 + \rho} \sqrt{\mathbb{E}_{P_{s,a}^*} (\eta - V(s'))_+^2} \right) \quad (19)$$

- Gradient w.r.t. η :

One trajectory

- [LMB⁺23] consider concrete cases: χ^2 and KL.
- For χ^2 case, optimize over λ :

$$\mathcal{T}_{\text{rob,c}}V = R(s, a) + \gamma \sup_{\eta} \left(\eta - \sqrt{1 + \rho} \sqrt{\mathbb{E}_{P_{s,a}^*} (\eta - V(s'))_+^2} \right) \quad (19)$$

- Gradient w.r.t. η :

$$\begin{aligned} g(\eta; V) &= 1 - \sqrt{1 + \rho} \frac{\mathbb{E}_{P_{s,a}^*} (\eta - V(s'))_+}{\sqrt{\mathbb{E}_{P_{s,a}^*} (\eta - V(s'))_+^2}} \\ &= 1 - \sqrt{1 + \rho} \frac{Z_1}{\sqrt{Z_2}} \end{aligned} \quad (20)$$

One trajectory

- Updating rule (at timestep t):

One trajectory

- Updating rule (at timestep t):

$$Z_{t+1,1}(s_t, a_t) = (1 - \alpha_{t,1})Z_{t,1}(s_t, a_t) + \alpha_{t,1}(\eta_t(s_t, a_t) - V_t(s_{t+1}))_+ \quad (21)$$

$$Z_{t+1,2}(s_t, a_t) = (1 - \alpha_{t,2})Z_{t,2}(s_t, a_t) + \alpha_{t,2}(\eta_t(s_t, a_t) - V_t(s_{t+1}))_+^2 \quad (22)$$

$$\eta_{t+1}(s_t, a_t) = (1 - \alpha_{t,3})\eta_t(s_t, a_t) + \alpha_{t,3}\left(1 - \sqrt{1 + \rho} \frac{Z_{t,1}(s_t, a_t)}{\sqrt{Z_{t,2}(s_t, a_t)}}\right) \quad (23)$$

$$Q_{t+1}(s_t, a_t) = (1 - \alpha_{t,4})Q_t(s_t, a_t) + \alpha_{t,4}(R(s_t, a_t) + \gamma(\eta_t(s_t, a_t) - \sqrt{1 + \rho}\sqrt{Z_{t,2}(s_t, a_t)})) \quad (24)$$

One trajectory

- Updating rule (at timestep t):

$$Z_{t+1,1}(s_t, a_t) = (1 - \alpha_{t,1})Z_{t,1}(s_t, a_t) + \alpha_{t,1}(\eta_t(s_t, a_t) - V_t(s_{t+1}))_+ \quad (21)$$

$$Z_{t+1,2}(s_t, a_t) = (1 - \alpha_{t,2})Z_{t,2}(s_t, a_t) + \alpha_{t,2}(\eta_t(s_t, a_t) - V_t(s_{t+1}))_+^2 \quad (22)$$

$$\eta_{t+1}(s_t, a_t) = (1 - \alpha_{t,3})\eta_t(s_t, a_t) + \alpha_{t,3}\left(1 - \sqrt{1 + \rho} \frac{Z_{t,1}(s_t, a_t)}{\sqrt{Z_{t,2}(s_t, a_t)}}\right) \quad (23)$$

$$Q_{t+1}(s_t, a_t) = (1 - \alpha_{t,4})Q_t(s_t, a_t) + \alpha_{t,4}(R(s_t, a_t) + \gamma(\eta_t(s_t, a_t) - \sqrt{1 + \rho}\sqrt{Z_{t,2}(s_t, a_t)})) \quad (24)$$

Theorem ([LMB⁺23])

As $t \rightarrow \infty$, $(Z_{t,1}, Z_{t,2}, \eta_t, Q_t)$ converges to $(Z_1^*, Z_2^*, \eta^*, Q^*)$ a.s.

One trajectory

- Assumption: $\alpha_{t,1} = \alpha_{t,2} = o(\alpha_{t,3})$, and $\alpha_{t,3} = o(\alpha_{t,4})$.

One trajectory

- Assumption: $\alpha_{t,1} = \alpha_{t,2} = o(\alpha_{t,3})$, and $\alpha_{t,3} = o(\alpha_{t,4})$.
- Updating Z_t with fixed η_t and Q_t :

One trajectory

- Assumption: $\alpha_{t,1} = \alpha_{t,2} = o(\alpha_{t,3})$, and $\alpha_{t,3} = o(\alpha_{t,4})$.
- Updating Z_t with fixed η_t and Q_t :

$$\begin{aligned}\dot{Z}(t) &= f(Z(t), \eta(t), Q(t)), \\ \dot{\eta}(t) &= 0, \dot{Q}(t) = 0.\end{aligned}\tag{25}$$

One trajectory

- Assumption: $\alpha_{t,1} = \alpha_{t,2} = o(\alpha_{t,3})$, and $\alpha_{t,3} = o(\alpha_{t,4})$.
- Updating Z_t with fixed η_t and Q_t :

$$\begin{aligned}\dot{Z}(t) &= f(Z(t), \eta(t), Q(t)), \\ \dot{\eta}(t) &= 0, \dot{Q}(t) = 0.\end{aligned}\tag{25}$$

$$Z(t) \rightarrow \lambda_1(\eta, Q).$$

One trajectory

- Assumption: $\alpha_{t,1} = \alpha_{t,2} = o(\alpha_{t,3})$, and $\alpha_{t,3} = o(\alpha_{t,4})$.
- Updating Z_t with fixed η_t and Q_t :

$$\begin{aligned}\dot{Z}(t) &= f(Z(t), \eta(t), Q(t)), \\ \dot{\eta}(t) &= 0, \dot{Q}(t) = 0.\end{aligned}\tag{25}$$

$Z(t) \rightarrow \lambda_1(\eta, Q)$.

- Updating η_t with fixed Q_t :

One trajectory

- Assumption: $\alpha_{t,1} = \alpha_{t,2} = o(\alpha_{t,3})$, and $\alpha_{t,3} = o(\alpha_{t,4})$.
- Updating Z_t with fixed η_t and Q_t :

$$\begin{aligned}\dot{Z}(t) &= f(Z(t), \eta(t), Q(t)), \\ \dot{\eta}(t) &= 0, \dot{Q}(t) = 0.\end{aligned}\tag{25}$$

$Z(t) \rightarrow \lambda_1(\eta, Q)$.

- Updating η_t with fixed Q_t :

$$\dot{\eta}(t) = g(\lambda_1(\eta(t), Q(t)), \eta(t), Q(t)), \dot{Q}(t) = 0.\tag{26}$$

One trajectory

- Assumption: $\alpha_{t,1} = \alpha_{t,2} = o(\alpha_{t,3})$, and $\alpha_{t,3} = o(\alpha_{t,4})$.
- Updating Z_t with fixed η_t and Q_t :

$$\begin{aligned}\dot{Z}(t) &= f(Z(t), \eta(t), Q(t)), \\ \dot{\eta}(t) &= 0, \dot{Q}(t) = 0.\end{aligned}\tag{25}$$

$$Z(t) \rightarrow \lambda_1(\eta, Q).$$

- Updating η_t with fixed Q_t :

$$\dot{\eta}(t) = g(\lambda_1(\eta(t), Q(t)), \eta(t), Q(t)), \dot{Q}(t) = 0.\tag{26}$$

$$\eta(t) \rightarrow \lambda_2(Q).$$

One trajectory

- Assumption: $\alpha_{t,1} = \alpha_{t,2} = o(\alpha_{t,3})$, and $\alpha_{t,3} = o(\alpha_{t,4})$.
- Updating Z_t with fixed η_t and Q_t :

$$\begin{aligned}\dot{Z}(t) &= f(Z(t), \eta(t), Q(t)), \\ \dot{\eta}(t) &= 0, \dot{Q}(t) = 0.\end{aligned}\tag{25}$$

$$Z(t) \rightarrow \lambda_1(\eta, Q).$$

- Updating η_t with fixed Q_t :

$$\dot{\eta}(t) = g(\lambda_1(\eta(t), Q(t)), \eta(t), Q(t)), \dot{Q}(t) = 0.\tag{26}$$

$$\eta(t) \rightarrow \lambda_2(Q).$$

- Finally,

$$\dot{Q}(t) = h(\lambda_1(Q(t)), \lambda_2(Q(t)), Q(t)).\tag{27}$$

- 1 Introduction
- 2 Model-free robust MDPs
- 3 Conclusion**
- 4 Reference

- 2 approaches to get a good estimator of non-linear objective:

- 2 approaches to get a good estimator of non-linear objective:
 - Multilevel Monte-Carlo

- 2 approaches to get a good estimator of non-linear objective:
 - Multilevel Monte-Carlo
 - Stochastic Gradient

- 2 approaches to get a good estimator of non-linear objective:
 - Multilevel Monte-Carlo
 - Stochastic Gradient
- Computation complexity for MMC is not considered.

- 2 approaches to get a good estimator of non-linear objective:
 - Multilevel Monte-Carlo
 - Stochastic Gradient
- Computation complexity for MMC is not considered.
- Stochastic gradient estimator needs a variant of robust MDPs, which still preserves robustness.

- 2 approaches to get a good estimator of non-linear objective:
 - Multilevel Monte-Carlo
 - Stochastic Gradient
- Computation complexity for MMC is not considered.
- Stochastic gradient estimator needs a variant of robust MDPs, which still preserves robustness.
- Data generating mechanism: from generative model to one trajectory. No finite-sample results for now.

- ① Introduction
- ② Model-free robust MDPs
- ③ Conclusion
- ④ Reference**

- [BGP19] Jose H Blanchet, Peter W Glynn, and Yanan Pei.
Unbiased multilevel monte carlo: Stochastic optimization, steady-state simulation, quantiles, and other applications.
arXiv preprint arXiv:1904.09929, 2019.
- [LBB⁺22] Zijian Liu, Qinxun Bai, Jose Blanchet, Perry Dong, Wei Xu, Zhengqing Zhou, and Zhengyuan Zhou.
Distributionally robust q -learning.
In *International Conference on Machine Learning*, pages 13623–13643. PMLR, 2022.
- [LMB⁺23] Zhipeng Liang, Xiaoteng Ma, Jose Blanchet, Jiheng Zhang, and Zhengyuan Zhou.
Single-trajectory distributionally robust reinforcement learning.
arXiv preprint arXiv:2301.11721, 2023.

- [LYZJ21] Xiang Li, Wenhao Yang, Zhihua Zhang, and Michael I Jordan.
Polyak-ruppert averaged q-leaning is statistically efficient.
arXiv preprint arXiv:2112.14582, 2021.
- [ND16] Hongseok Namkoong and John C Duchi.
Stochastic gradient methods for distributionally robust optimization with f-divergences.
Advances in neural information processing systems, 29, 2016.
- [Wai19] Martin J Wainwright.
Stochastic approximation with cone-contractive operators: Sharp l_∞ -bounds for q -learning.
arXiv preprint arXiv:1905.06265, 2019.

- [WSBZ23] Shengbo Wang, Nian Si, Jose Blanchet, and Zhengyuan Zhou.
A finite sample complexity bound for distributionally robust q-learning.
In *International Conference on Artificial Intelligence and Statistics*, pages 3370–3398. PMLR, 2023.
- [YWK⁺23] Wenhao Yang, Han Wang, Tadashi Kozuno, Scott M Jordan, and Zhihua Zhang.
Avoiding model estimation in robust markov decision processes with a generative model.
arXiv preprint arXiv:2302.01248, 2023.
- [YZZ22] Wenhao Yang, Liangyu Zhang, and Zhihua Zhang.
Toward theoretical understandings of robust Markov decision processes: Sample complexity and asymptotics.
The Annals of Statistics, 50(6):3223 – 3248, 2022.

