

Robust Offline Reinforcement Learning via Lipschitz Value Function

Tengyang Xie

tx10@illinois.edu

Overview

Goal: Provably Robust RL against (adversarial) perturbations on state observations.

Results:

1. On robustness definition: Using Lipschitz Q-function is sufficient against perturbations on states (done: formally proved)
2. Lipschitz policy evaluation
(method proposed with basic guarantee. future work: detailed analysis)
3. Robust offline RL (Lipschitz policy evaluation + Policy improvement) (future work)

Lipschitz Value Function

Fact: Lipschitz (transition + reward) \implies Lipschitz Q-Function
[Asadi et al., 2018]

Theorem. [Policy improvement via Lipschitz Q-function is robust against perturbations on states]

Assume (\hat{Q}, π) pair satisfies (1) $\|\hat{Q} - Q^\pi\|_\infty = \varepsilon_1$; (2) \hat{Q} is L -Lipschitz. Let $\pi_{\hat{Q}} \circ \nu_{\varepsilon_2}$ be the perturbed greedy policy w.r.t. \hat{Q} and an ε_2 -bounded perturbation ν_{ε_2} , then

$$V^{\pi_{\hat{Q}} \circ \nu_{\varepsilon_2}}(s) \geq V^\pi(s) - \frac{2\varepsilon_1 + 2L\varepsilon_2}{1 - \gamma}, \quad \forall s \in \mathcal{S},$$

where $\pi \circ \nu_{\varepsilon_2}(s) = \operatorname{argmax}_{a \in \mathcal{A}} \hat{Q}(\nu_{\varepsilon_2}(s), a)$.

Lipschitz Value Function

What can we learn from that Theorem?

Policy improvement still holds approximately even with adversarial perturbations on states as long as:

- (i) We can approximate Q^π accurately;
- (ii) Approximated Q^π is Lipschitz continuous.

A robust RL framework based on that concept:

1. Policy evaluation with ensuring Lipschitz; (next page)
2. Policy improvement according to the estimated Lipschitz Q-function. (future work)

Lipschitz Policy Evaluation

A non-parametric Lipschitz policy evaluation approach:

Given batch data $\mathcal{D} = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^n$ (all the states are perturbed by v_{ε_2}) and target policy π .

1. Set $q_{0,i} = r_i - \gamma L d(s_i, s'_i), \forall i \in [n]$
2. For $t = 1, 2, \dots, T$
 - 2.1 $Q_t(s, a) = \max_{i \in [n]: a_i = a} (q_{t-1,i} - L d(s, s_i) - L \varepsilon_2), \forall (s, a) \in \mathcal{S} \times \mathcal{A}$
 - 2.2 $q_{t,i} = r_i + \gamma [\hat{T}^\pi Q_t](s_i, a_i), \forall i \in [n]$

Theorem

Let $Q = Q_T, T \rightarrow \infty$, we have

$$\|Q_T - Q^\pi\|_\infty \leq \frac{2L(\varepsilon_S + \varepsilon_2)}{1 - \gamma},$$

where $\varepsilon_S := \sup_{(s,a) \in \mathcal{S} \times \mathcal{A}} \min_{i \in [n]: a_i = a} d(s, s_i)$, and ε_2 is the perturbation bound.

Reference

Kavosh Asadi, Dipendra Misra, and Michael Littman. Lipschitz continuity in model-based reinforcement learning. In *International Conference on Machine Learning*, pages 264–273, 2018.