

---

# RORL: Robust Offline Reinforcement Learning via Conservative Smoothing

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Offline reinforcement learning (RL) provides a promising direction to exploit the  
2 massive amount of offline data for complex decision-making tasks. Due to the  
3 distribution shift issue, current offline RL algorithms are generally designed to be  
4 conservative for value estimation and action selection. However, such conservatism  
5 impairs the robustness of learned policies, leading to a significant change even for  
6 a small perturbation on observations. To trade off robustness and conservatism, we  
7 propose Robust Offline Reinforcement Learning (RORL) with a novel conservative  
8 smoothing technique. In RORL, we explicitly introduce regularization on the policy  
9 and the value function for states near the dataset and additional conservative value  
10 estimation on these OOD states. Theoretically, we show RORL enjoys a tighter  
11 suboptimality bound than recent theoretical results in linear MDPs. We demonstrate  
12 that RORL can achieve the state-of-the-art performance on the general offline RL  
13 benchmark and is considerably robust to adversarial observation perturbation.

## 14 1 Introduction

15 Over the past few years, deep reinforcement learning (RL) has been a vital tool for various decision-  
16 making tasks [32, 45, 43, 9] in a trial-and-error manner. A major limitation of current deep RL  
17 algorithms is that they require intense online interactions with the environment [28]. These data  
18 collecting processes can be costly and even prohibitive in many real-world scenarios such as robotics  
19 and health care. Offline RL [13, 26] is gaining more attention recently since it offers probabilities to  
20 learn reinforced decision-making strategies from fully offline datasets.

21 The main challenge of offline RL is the distribution shift between the offline dataset and the learned  
22 policy, which would lead to significant overestimation for the out-of-distribution (OOD) actions  
23 [13, 26]. To overcome such an issue, a series of offline RL works [52, 13, 60, 27, 2, 4] propose to  
24 celebrate conservatism, such as constraining the learned policy close to supported distribution or  
25 penalizing the  $Q$ -values of OOD actions. In addition, another stream of works builds upon model-  
26 based offline RL algorithms [63, 62, 51], which leverages the ensemble dynamics models to enforce  
27 pessimism via uncertainty penalizing or data generation.

28 However, conservatism is not the only concern when applying offline RL to the real world. Due  
29 to the sensor errors and model mismatch, the robustness of offline RL is crucial under the realistic  
30 engineering conditions, which has not been well studied yet. In online RL, a series of works have  
31 been studied to learn the optimal policy under worst-case perturbations of the observation [64, 37, 18]  
32 or environmental dynamics [50, 39, 40]. Yet, it is non-trivial to apply online robust RL techniques  
33 into the offline problems. The main challenge is that the perturbation of either states or actions may  
34 bring OOD data and extra overestimation for the value function. New techniques are needed to tackle  
35 the conservatism and robustness simultaneously in the offline RL.

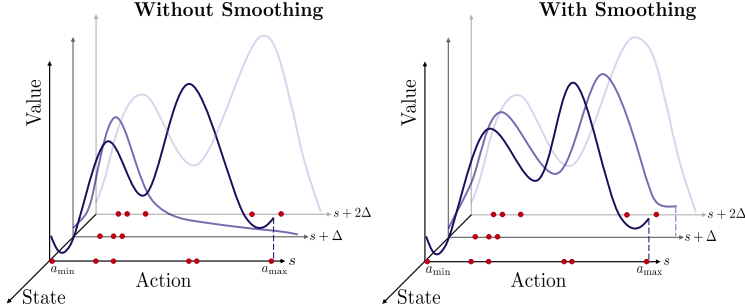


Figure 1: A schematic diagram of RORL. The red spots represent the offline data samples. Without state smoothing, the value function would change drastically over neighbor states and induce an unstable policy. Yet, the smoothness may also lead to value overestimation of dangerous areas. RORL trades off smoothness and possible overestimation as discussed in Sec 4.

36 This paper studies robust offline RL against adversarial observation perturbation, where the agent  
 37 needs to learn the policy conservatively while handling the potential OOD observation with per-  
 38 turbation. We first demonstrate that current value-based offline RL algorithms lack the necessary  
 39 smoothness for the policy, which is visualized in Figure 1. As an illustration, we show that a famous  
 40 baseline method CQL [27] learns a non-smooth value function, leading to significant performance  
 41 degradation for even a tiny scale perturbation on observation (see Section 3 for details). Notice that  
 42 simply adopting the smoothing technique for existing methods may result in extra overestimation at  
 43 the boundary of supported distribution and lead the agent toward the unsafe areas.

44 To this end, we propose Robust Offline Reinforcement Learning (RORL) with a novel conservative  
 45 smoothing technique, which explicitly handles the overestimation of OOD state-action pairs. Specifi-  
 46 cally, we explicitly introduce smooth regularization on both the values and policies for states near  
 47 the dataset support and conservatively estimate the values of these OOD states based on pessimistic  
 48 bootstrapping. In addition, we theoretically prove that RORL yields a valid uncertainty quantifier in  
 49 linear MDPs and enjoys a tighter suboptimality bound than previous work [4].

50 In our experiments, we demonstrate that RORL can achieve state-of-the-art (SOTA) performance  
 51 in D4RL benchmark [10] for offline RL given a suitable adversarial perturbation scale during  
 52 training. The surprising performance in the benchmark shows that the robust training improves the  
 53 generalization ability and leads to performance improvement even in non-perturbed tasks. Meanwhile,  
 54 compared with the current SOTA baseline [2], RORL is considerably more robust to adversarial  
 55 perturbations on observations. We conduct the ablation study under different types of adversarial  
 56 attacks on the observations, showing the consistently superior performance over several tasks.

## 57 2 Preliminaries

58 **Offline RL and PBRL** Considering an episodic MDP  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, T, r, \gamma, \mathbb{P})$ , where  $\mathcal{S}$  is the state  
 59 space,  $\mathcal{A}$  is the action space,  $T$  is the length of the episode,  $r$  is the reward function,  $\mathbb{P}$  is the dynamics,  
 60 and  $\gamma$  is the discount factor. In offline RL, the objective of the agent is to find an optimal policy by  
 61 sampling experience from a fixed dataset  $\mathcal{D} = \{(s_t^i, a_t^i, r_t^i, s_{t+1}^i)\}$ . Nevertheless, directly applying  
 62 off-policy algorithms in offline RL has distribution shift problem. In  $Q$ -learning, the value function  
 63 evaluated on the greedy action  $a'$  in Bellman operator  $\mathcal{T}Q = r + \gamma \mathbb{E}_{s'}[\max_{a'}(s', a')]$  tends to have  
 64 extrapolation error since  $(s', a')$  has barely occurred in  $\mathcal{D}$ . There are mainly three kinds of methods  
 65 to track the distribution shift problem in offline RL, including policy constraints, conservative value  
 66 function, and uncertainty-based methods.

67 Pessimistic Bootstrapping for Offline RL (PBRL) [4] is an uncertainty-based method that uses boot-  
 68 strapped  $Q$ -functions for uncertainty quantification and OOD sampling for regularization. Specifically,  
 69 PBRL maintains  $K$  bootstrapped  $Q$  functions to quantify the epistemic uncertainty and performs  
 70 pessimistic update to penalize  $Q$  functions with large uncertainties. The uncertainty is defined as  
 71 the standard deviation among bootstrapped  $Q$ -functions. For each bootstrapped  $Q$ -function, the  
 72 Bellman target is defined as  $\hat{\mathcal{T}}Q(s, a) = r(s, a) + \gamma \hat{\mathbb{E}}_{s' \sim P(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q(s', a') - \beta u(s', a')]$ .  
 73 Under linear MDP assumptions, this uncertainty is equivalent to the LCB penalty and is provably  
 74 efficient [22]. Furthermore, PBRL incorporates OOD sampling in training by sampling OOD actions

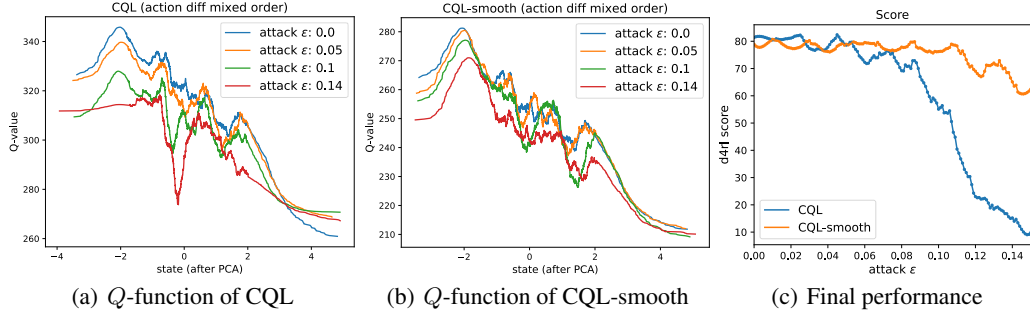


Figure 2: Figure (a) and (b) illustrate the  $Q$ -functions of  $\hat{s}$  with adversarial noises in CQL and CQL-smooth, respectively. The same moving average factor is used in plotting both figures. Figure (c) shows the performance evaluation of CQL and CQL-smooth with different perturbation scale. We use 100 different  $\epsilon \in [0.0, 0.15]$  for the evaluation.

75 to form  $(s, a^{\text{ood}})$  pairs, where  $a^{\text{ood}}$  follows the learned policy. The learning target for  $(s, a^{\text{ood}})$  is  
 76  $\hat{\mathcal{T}}^{\text{ood}}Q(s, a^{\text{ood}}) := Q(s, a^{\text{ood}}) - \beta u(s, a^{\text{ood}})$ , which introduces uncertainty penalization to enforce  
 77 pessimistic  $Q$ -functions for OOD actions.

78 **Smooth Regularized RL** Robust RL aims to learn a robust policy against the adversarial per-  
 79 turbed environment in online RL. SR<sup>2</sup>L [44] enforces smoothness in both the policy and  $Q$ -  
 80 functions. Specifically, SR<sup>2</sup>L encourages the outputs of the policy and value function to not  
 81 change much when injecting small perturbation to the state. For state  $s$ , SR<sup>2</sup>L constructs a per-  
 82 turbation set  $\mathbb{B}_d(s, \epsilon) = \{\hat{s} : d(s, \hat{s}) \leq \epsilon\}$  and introduces a smoothness regularizer for policy as  
 83  $\mathcal{R}_s^\pi = \mathbb{E}_{s \sim \rho^\pi} \max_{\hat{s} \in \mathbb{B}_d(s, \epsilon)} \mathcal{D}(\pi(\cdot|s) \parallel \pi(\cdot|\hat{s}))$ , where  $\mathcal{D}(\cdot|\cdot)$  is a distance metric and the max oper-  
 84 ator gives an adversarial manner to choose  $\hat{s}$ . Similarly, the smoothness regularizer for the value  
 85 function is defined as  $\mathcal{R}_s^V = \mathbb{E}_{s \sim \rho^\pi, a \sim \pi} \max_{\hat{s} \in \mathbb{B}_d(s, \epsilon)} (Q(s, a) - Q(\hat{s}, a))^2$ . The induced smoothness  
 86 is shown to improve robustness against both random and adversarial perturbations.

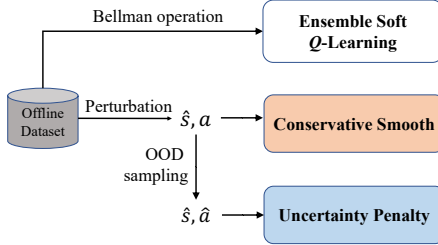
### 87 3 Robustness of Offline RL: A Motivating Example

88 We give a motivating example to illustrate the robustness of the popular CQL [27] policies. We  
 89 introduce an adversarial attack on state  $s$  to obtain  $\hat{s} = \arg \max_{\hat{s} \in \mathbb{B}_d(s, \epsilon)} D_J(\pi_\theta(\cdot|s) \parallel \pi_\theta(\cdot|\hat{s}))$ , where  
 90  $\mathbb{B}_d$  is the perturbation set to ensure  $d(s, \hat{s}) \leq \epsilon$  and the Jeffrey’s divergence  $D_J$  for two distributions  
 91  $P, Q$  is defined by:  $D_J(P \parallel Q) = \frac{1}{2}[D_{\text{KL}}(P \parallel Q) + D_{\text{KL}}(Q \parallel P)]$ . To obtain  $\hat{s}$ , we take gradient  
 92 ascent with respect to the loss function  $D_J(\pi_\theta(\cdot|s) \parallel \pi_\theta(\cdot|\hat{s}))$  and restrict the outputs to the  $\mathbb{B}_d(s, \epsilon)$   
 93 set, where  $\pi_\theta$  is a learned CQL policy and  $\epsilon$  controls the strength of attack.

94 In the *walker-medium-v2* task from D4RL [10], we use various setups of  $\epsilon$  for adversarial attack to  
 95 evaluate the robustness of CQL policies. Specifically, we use  $\epsilon \in \{0, 0.05, 0.1, 0.14\}$  to control the  
 96 strengths of the attack, where we have  $\hat{s} = s$  if  $\epsilon = 0$ . Given a specific  $\epsilon$ , we sample  $N$  state-action  
 97 pairs  $\{(s_i, a_i)\}$  from the offline dataset, and then perform adversarial attack to obtain  $\{(\hat{s}_i, a_i)\}$  and  
 98 also the corresponding  $Q$ -values  $\{Q_i(\hat{s}_i, a_i)\}$ , where the  $Q$ -function is the trained CQL critic.

99 Figure 2(a) shows the relationship between  $\hat{s}_i$  and the corresponding  $Q_i$  with different  $\epsilon$ . To visualize  
 100  $\hat{s}_i$ , we perform PCA dimensional reduction [48] and choose one of the reduced dimensions to  
 101 represent  $\hat{s}_i$ . With the increase of  $\epsilon$  in the adversarial attack, the  $Q$ -curve has greater deviation  
 102 compared to the curve with  $\epsilon = 0$ . The result signifies that the  $Q$ -function of CQL is not smooth in  
 103 state space, which makes the adversarial noises easily affect the  $Q$ -function. As a comparison, we  
 104 apply the proposed conservative smoothing loss in CQL training (i.e., *CQL-smooth*) and use the same  
 105 evaluation method to obtain  $\hat{s}_i$  and  $Q_i$ . According to the result in Figure 2(b), the value function  
 106 becomes smoother.

107 In addition, we show how the adversarial attack affects the final performance of offline RL policies.  
 108 We use  $\epsilon \in [0, 0.15]$  to evaluate both the original CQL policies (i.e., *CQL*) and CQL with conservative  
 109 smoothing loss (i.e., *CQL-smooth*) in adversarial attack. Figure 2(c) shows the performance with  
 110 different settings of  $\epsilon$ . We find that our smooth constraints significantly improve the robustness of  
 111 CQL, especially for large adversarial noises.




---

**Algorithm 1: RORL**


---

**while not converged do**

- Sample mini-batch transitions  $(s, a, r, s')$  from  $\mathcal{D}$ .
  - Perform extended soft  $Q$ -learning to train multiple  $Q$ -functions  $\{Q_1, \dots, Q_K\}$ .
  - Sample  $\hat{s}$  from  $\mathbb{B}_d(s, \epsilon)$  to obtain  $(\hat{s}, a)$  pairs.
  - Perform adversarial smoothing for the  $Q$ -network.
  - Calculate uncertainty to penalize  $(\hat{s}, \hat{a})$  for conservatism.
  - Combining terms to train the critic with Eq. (5).
  - Train the policy with smooth constraint as Eq. (6).
- 

Figure 3: **RORL Algorithm:** RORL trains multiple soft  $Q$ -functions for uncertainty quantification. The conservative smoothing loss is calculated for  $(\hat{s}, a)$  with perturbed state. We perform uncertainty penalization for  $(\hat{s}, \hat{a})$  with OOD actions.

## 122 4 Robust Offline RL via Conservative Smoothing

123 In RORL, we develop smooth regularization on both the policy and the value function for states  
 124 near the dataset. The smooth constraints make the policy and  $Q$ -function robust to perturbation in  
 125 the dataset (for offline training) and interactive samples (for online evaluation). Nevertheless, the  
 126 smoothness may lead to value overestimation of the edge of the supported dataset. To address this  
 127 problem, we adopt bootstrapped  $Q$ -functions [35] for uncertainty quantification and sample OOD  
 128 actions for penalization. RORL obtains conservative and smooth value estimation on both OOD  
 129 states and actions, which improves the generalization ability of offline RL algorithms. The overall  
 130 architecture of RORL is given in Figure 3.

131 **Robust  $Q$ -function** We sample three sets of state-action pairs and apply different loss functions  
 132 to obtain a conservative smooth policy. Specifically, for a  $(s, a)$  pair sampled from  $\mathcal{D}$ , we construct  
 133 a perturbation set  $\mathbb{B}_d(s, \epsilon)$  to obtain  $(\hat{s}, a)$  pairs, where  $\hat{s} \in \mathbb{B}_d(s, \epsilon)$  and  $\epsilon$ . Then we perform OOD  
 134 sampling by using the current policy  $\pi_\theta$  to obtain  $(\hat{s}, \hat{a})$  pairs, where  $\hat{a} \sim \pi_\theta(\hat{s})$ . RORL contains  $K$   
 135  $Q$ -functions trained by Soft  $Q$ -learning as critics. We denote the parameters of  $i$ -th  $Q$ -function and  
 136 the target  $Q$ -function as  $\phi_i$  and  $\phi'_i$ , respectively. In the following, we give different learning targets  
 137 for  $(s, a)$ ,  $(\hat{s}, a)$ , and  $(\hat{s}, \hat{a})$  pairs, respectively.

138 First, for a  $(s, a)$  pair sampled from  $\mathcal{D}$ , we apply extended soft  $Q$ -learning to obtain the target as

$$\widehat{T}Q_{\phi_i}(s, a) := r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_\theta(\cdot|s)} \left[ \min_{j=1, \dots, K} Q_{\phi'_j}(s', a') - c \cdot \log \pi_\theta(a'|s') \right], \quad (1)$$

139 where the next- $Q$  function takes minimum value among the target  $Q$ -functions.

140 Then, for a  $(\hat{s}, a)$  pair with a perturbed state, we enforce smoothness in  $Q$ -function by minimizing the  
 141  $Q$ -value difference between  $Q(s, a)$  and  $Q(\hat{s}, a)$ . In particular, we choose an adversarial  $\hat{s} \in \mathbb{B}_d(s, \epsilon)$   
 142 that maximizes  $\mathcal{L}(Q(\hat{s}, a), Q(s, a))$ , and then minimize this loss function with the adversarial  $\hat{s}$ ,  
 143 where  $\mathcal{L}$  is a distance measure. Intuitively, we want that the  $Q$ -function is smooth with the most  
 144 difficult (i.e., adversarial) attack in  $\mathbb{B}_d(s, \epsilon)$ , where  $\epsilon$  controls the perturbation scale. To solve such a  
 min-max problem, we maximize  $\mathcal{L}(Q(\hat{s}, a), Q(s, a))$  to obtain  $\hat{s}$ , and then minimize the loss function  
 with the selected  $\hat{s}$ .

$$\mathcal{L}_{\text{smooth}}(s, a; \phi_i) = \max_{\hat{s} \in \mathbb{B}_d(s, \epsilon)} \mathcal{L}(Q_{\phi_i}(\hat{s}, a), Q_{\phi_i}(s, a)) \quad (2)$$

145 where  $\mathcal{L}(\cdot, \cdot)$  is the smooth loss function, and we denote  $\delta(s, \hat{s}, a) = Q_{\phi_i}(\hat{s}, a) - Q_{\phi_i}(s, a)$ . We  
 146 remark that if  $\delta(s, \hat{s}, a) > 0$ , the perturbed state may induce an overestimated  $Q$ -value that we need  
 147 to smooth. In contrast, if  $\delta(s, \hat{s}, a) < 0$ , the perturbed  $Q$ -function is underestimated, which does  
 148 not cause a serious problem in offline RL. As a result, we use different weights for  $\delta(s, \hat{s}, a)_+$  and  
 149  $\delta(s, \hat{s}, a)_-$ . The definition of  $\mathcal{L}(\cdot, \cdot)$  is give as follows,

$$\mathcal{L}(Q_{\phi_i}(\hat{s}, a), Q_{\phi_i}(s, a)) = (1 - \tau)\delta(s, \hat{s}, a)_+^2 + \tau\delta(s, \hat{s}, a)_-^2, \quad (3)$$

150 where we choose  $\tau \leq 0.5$ ,  $x_+ = \max(x, 0)$  and  $x_- = \min(x, 0)$ . In this term, we does not introduce  
 151 OOD action  $\hat{a}$  for smoothing since the actions are desired to be close to the behavior actions in the  
 152 dataset for areas near the offline data.

145 Finally, to prevent overestimation of OOD states and actions, we use bootstrapped uncertainty  $u(\hat{s}, \hat{a})$   
 146 as a penalty for  $Q(\hat{s}, \hat{a})$ , where  $\hat{a} \sim \pi(\hat{s})$  is an OOD action sampled from the current policy  $\pi$ .  
 147 We remark that a similar OOD sampling is also used in PBRL [4]. *The difference is that PBRL*  
 148 *only penalizes the OOD actions, while RORL penalizes both the OOD states and actions to provide*  
 149 *conservatism for unfamiliar areas.* The ensemble technique in RL comes from Bootstrapped DQN  
 150 [35]. The ensemble forms an estimation of the  $Q$ -posterior, which yields diverse predictions on areas  
 151 with scarce data.  $u(\hat{s}, \hat{a})$  quantifies the deviation of a datapoint from the offline dataset [2, 4]. We  
 152 follow PBRL and use a loss function as:

$$\mathcal{L}_{\text{ood}}(s; \phi_i) = \mathbb{E}_{\hat{s} \sim \mathbb{B}_d(s, \epsilon), \hat{a} \sim \pi_\theta(\hat{s})} (\widehat{\mathcal{T}}_{\text{ood}} Q_{\phi_i}(\hat{s}, \hat{a}) - Q_{\phi_i}(\hat{s}, \hat{a}))^2, \quad (4)$$

153 where the pseudo-target for the OOD datapoints is computed as:  $\widehat{\mathcal{T}}_{\text{ood}} Q_{\phi_i}(\hat{s}, \hat{a}) := Q_{\phi_i}(\hat{s}, \hat{a}) -$   
 154  $u(\hat{s}, \hat{a})$ , where  $u(\hat{s}, \hat{a}) := \sqrt{\frac{1}{K} \sum_{k=1}^K (Q_{\phi_k}(\hat{s}, \hat{a}) - \bar{Q}_\phi(\hat{s}, \hat{a}))^2}$ , The bootstrapped uncertainty  
 155  $u(\hat{s}, \hat{a})$  is defined as the standard deviation among the  $Q$ -ensemble.

156 Combining the loss functions above, RORL has the following loss function for each  $Q_{\phi_i}$ :

$$\min_{\phi_i} \mathbb{E}_{s, a, r, s' \sim \mathcal{D}} \left[ (\widehat{\mathcal{T}} Q_{\phi_i}(s, a) - Q_{\phi_i}(s, a))^2 + \alpha \mathcal{L}_{\text{smooth}}(s, a; \phi_i) + \beta \mathcal{L}_{\text{ood}}(s; \phi_i) \right], \quad (5)$$

157 **Robust Policy** We learn a robust policy by using a smooth constraint to make the policy change  
 158 less with perturbations. Similarly, we choose an adversarial state  $\hat{s} \in \mathbb{B}_d(s, \epsilon)$  that maximizes  
 159  $D_J(\pi_\theta(\cdot|s) \parallel \pi_\theta(\cdot|\hat{s}))$ , and then minimize the policy difference between  $\pi_\theta(\cdot|s)$  and  $\pi_\theta(\cdot|\hat{s})$ . To  
 160 conclude, we minimize the following loss function for the policy,

$$\min_{\theta} \left[ \mathbb{E}_{s, a, r, s' \sim \mathcal{D}} \left[ - \min_{j=1, \dots, K} Q_{\phi_j}(s, a) + \alpha_2 \max_{\hat{s} \in \mathbb{B}_d(s, \epsilon)} D_J(\pi_\theta(\cdot|s) \parallel \pi_\theta(\cdot|\hat{s})) + c \log \pi_\theta(a|s) \right] \right]. \quad (6)$$

161 where the first term aims to maximize the minimum of the ensemble  $Q$ -functions to obtain a  
 162 conservative policy, and the third term is regularization of entropy.

## 163 5 Theoretical Analysis

164 We analyze the simplified learning objective of RORL in linear MDPs [21, 22], where the feature  
 165 map of the state-action pair takes the form of  $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ , and the transition kernel and the  
 166 reward function are assumed to be linear in  $\phi$ . The parameter  $\tilde{w}_t$  of RORL can be solved in closed  
 167 form following the least squares value iteration algorithm (LSVI), which minimizes the following  
 168 loss function.

$$\begin{aligned} \tilde{w}_t^i = \min_{w \in \mathcal{R}^d} & \left[ \sum_{i=1}^m (y_t^i - Q_w(s_t^i, a_t^i))^2 + \sum_{i=1}^m \frac{1}{|\mathbb{B}_d(s_t^i, \epsilon)|} \sum_{\hat{s}_t^i \in \mathcal{D}_{\text{ood}}(s_t^i)} (Q_w(s_t^i, a_t^i) - Q_w(\hat{s}_t^i, a_t^i))^2 + \right. \\ & \left. \sum_{(\hat{s}, \hat{a}, \hat{y}) \sim \mathcal{D}_{\text{ood}}} (\hat{y} - Q_w(\hat{s}, \hat{a}))^2 \right], \end{aligned} \quad (7)$$

169 where we have  $Q_w(s_t^i, a_t^i) = \phi(s_t^i, a_t^i)^\top w$  since the  $Q$ -function is also linear in  $\phi$ . The first term  
 170 in Eq. (7) is the ordinary TD-error, where we consider the setting of  $\gamma = 1$  and the  $Q$ -target  
 171 is  $y_t^i = r(s_t^i, a_t^i) + V_{t+1}(s_{t+1}^i)$ . The second term is the proposed conservative smoothing loss.  
 172 Specifically,  $\hat{s}_t^i \sim \mathcal{D}_{\text{ood}}(s_t^i)$  are sampled from a  $l_\infty$  ball of center  $s_t^i$  and norm  $\epsilon > 0$ , which can be  
 173 formulated as  $\hat{s} \sim \mathbb{B}_d(s_t^i, \epsilon)$ . The third term is the additional OOD-sampling loss [4], which enforces  
 174 conservatism for OOD actions. Differently, we use perturbed states sampled from  $\mathcal{D}_{\text{ood}} = \bigcup_{i=1}^m \mathcal{D}_{\text{ood}}(s_t^i)$   
 175 rather than states from dataset in PBRL. The OOD action  $\hat{a}$  is sampled from policy  $\pi$ . The explicit  
 176 solution of Eq. (7) takes the following form by following LSVI:

$$\tilde{w}_t^i = \tilde{\Lambda}_t^{-1} \left( \sum_{i=1}^m \phi(s_t^i, a_t^i) y_t^i + \sum_{(\hat{s}, \hat{a}, \hat{y}) \sim \mathcal{D}_{\text{ood}}} \phi(\hat{s}, \hat{a}) \hat{y} \right), \quad (8)$$

177 where the covariance matrix  $\tilde{\Lambda}_t$  is defined as

$$\begin{aligned} \tilde{\Lambda}_t = & \sum_{i=1}^m \phi(s_t^i, a_t^i) \phi(s_t^i, a_t^i)^\top + \sum_{(\hat{s}, \hat{a}) \sim \mathcal{D}_{\text{ood}}} \phi(\hat{s}_t, \hat{a}_t) \phi(\hat{s}_t, \hat{a}_t)^\top \\ & + \sum_{i=1}^m \frac{1}{|\mathbb{B}_d(s_t^i, \epsilon)|} \sum_{\hat{s}_t^i \sim \mathcal{D}_{\text{ood}}(s_t^i)} [\phi(\hat{s}_t^i, a_t^i) - \phi(s_t^i, a_t^i)] [\phi(\hat{s}_t^i, a_t^i) - \phi(s_t^i, a_t^i)]^\top. \end{aligned} \quad (9)$$

178 We denote the first term and the second term as  $\tilde{\Lambda}_t^{\text{in}}$  and  $\tilde{\Lambda}_t^{\text{ood}}$ , which represent the covariance matrices  
 179 induced by the offline samples and OOD samples, respectively. Nevertheless, in linear MDPs, it is  
 180 difficult to ensure the covariance  $\tilde{\Lambda}_t^{\text{in}} + \tilde{\Lambda}_t^{\text{ood}} \succeq \lambda I$ , since it requires that the embeddings of the samples  
 181 are isotropic to make the eigenvalues of the corresponding covariance matrix lower bounded. This  
 182 condition holds if we can sample embeddings uniformly from the whole embedding space. However,  
 183 since the offline dataset has limited coverage in the state-action space and the OOD samples come  
 184 from a limited  $l_\infty$ -ball around the offline data,  $\tilde{\Lambda}_t^{\text{in}} + \tilde{\Lambda}_t^{\text{ood}}$  cannot be guaranteed to be positive definite.  
 185 PBRL [4] uses the assumption of  $\tilde{\Lambda}_t^{\text{ood}} \succeq \lambda I$ , while it is unachievable empirically. In RORL, we solve  
 186 this problem by introducing an additional conservative smoothing loss, which induces a covariance  
 187 matrix as  $\tilde{\Lambda}_t^{\text{ood-diff}} = \sum_{i=1}^m \frac{1}{|\mathbb{B}_d(s_t^i, \epsilon)|} \sum_{\hat{s}_t^i \sim \mathcal{D}_{\text{ood}}(s_t^i)} [\phi(\hat{s}_t^i, a_t^i) - \phi(s_t^i, a_t^i)] [\phi(\hat{s}_t^i, a_t^i) - \phi(s_t^i, a_t^i)]^\top$  (i.e.,  
 188 the third term in Eq. (9)). The following theorem gives the guarantees of  $\tilde{\Lambda}_t^{\text{ood-diff}} \succeq \lambda I$ .

189 **Theorem 1.** Assume  $\exists i \in [1, m]$  the vector group of all  $\hat{s}_t^i \sim \mathcal{D}_{\text{ood}}(s_t^i)$ :  $\{\phi(\hat{s}_t^i, a_t^i) - \phi(s_t^i, a_t^i)\}$  be  
 190 full rank, then the covariance matrix  $\tilde{\Lambda}_t^{\text{ood-diff}}$  is positive-definite:  $\tilde{\Lambda}_t^{\text{ood-diff}} \succeq \lambda \cdot I$  where  $\lambda > 0$ .

191 Recall the covariance matrix of PBRL is  $\tilde{\Lambda}_t^{\text{PBRL}} = \tilde{\Lambda}_t^{\text{in}} + \tilde{\Lambda}_t^{\text{ood}}$ , and RORL has a covariance matrix as  
 192  $\tilde{\Lambda}_t = \tilde{\Lambda}_t^{\text{PBRL}} + \tilde{\Lambda}_t^{\text{ood-diff}}$ , we have the following corollary based on Theorem 1

193 **Corollary 1.** Under the linear MDP assumptions, we have  $\tilde{\Lambda}_t \succeq \tilde{\Lambda}_t^{\text{PBRL}}$ . Further, the covariance  
 194 matrix  $\tilde{\Lambda}_t$  of RORL is positive-definite:  $\tilde{\Lambda}_t \succeq \lambda \cdot I$ , where  $\lambda > 0$ .

195 Recent theoretical analysis shows that an appropriate uncertainty quantification is essential to provable  
 196 efficiency in offline RL [22, 58, 4]. Pessimistic Value Iteration [22] defines a general  $\xi$ -uncertainty  
 197 quantifier as the penalty and achieves provable efficient pessimism in offline RL. In linear MDPs,  
 198 Lower Confidence Bound (LCB)-penalty [1, 21] is known to be a  $\xi$ -uncertainty quantifier for  
 199 appropriately selected  $\beta_t$  as  $\Gamma^{\text{lcb}}(s_t, a_t) = \beta_t \cdot [\phi(s_t, a_t)^\top \Lambda_t^{-1} \phi(s_t, a_t)]^{1/2}$ . Following the analysis  
 200 of PBRL [4], since the bootstrapped uncertainty is an estimation to the LCB-penalty and the OOD  
 201 sampling provides a covariance matrix of  $\tilde{\Lambda}_t^{\text{PBRL}}$ , the proposed RORL also form a valid  $\xi$ -uncertainty  
 202 quantifier with the covariance matrix  $\tilde{\Lambda}_t \succeq \lambda \cdot I$  given in Corollary 1. This allows us to further  
 203 characterize the optimality gap based on the pessimistic value iteration [22, 4]. We have the following  
 204 suboptimality gap under linear MDP assumptions.

205 **Corollary 2.**  $\text{SubOpt}(\pi^*, \hat{\pi}) \leq \sum_{t=1}^T \mathbb{E}_{\pi^*} [\Gamma_t^{\text{lcb}}(s_t, a_t)] < \sum_{t=1}^T \mathbb{E}_{\pi^*} [\Gamma_t^{\text{lcb-PBRL}}(s_t, a_t)]$ .

206 We refer to Appendix A for detailed proof. The second inequality is directly induced by  $\Gamma_t^{\text{lcb}}(s_t, a_t) <$   
 207  $\Gamma_t^{\text{lcb-PBRL}}(s_t, a_t)$ . Therefore, RORL enjoys a tighter suboptimality bound than PBRL [4].

## 208 6 Experiments

209 We evaluate our method on D4RL benchmark [11] with various continuous-control tasks and datasets.  
 210 We compare RORL with several offline RL algorithms, including (i) BC that performs behavior  
 211 cloning, (ii) CQL [27] that learns conservative value function for OOD actions, (iii) EDAC [2] that  
 212 learns a diversified  $Q$ -ensemble to enforce conservatism, and (iv) PBRL [4] that performs uncertainty  
 213 penalization and OOD sampling. Among these methods, EDAC [2] and PBRL [4] are related to  
 214 RORL since all these methods apply  $Q$ -ensemble for conservatism. EDAC needs much more  $Q$ -  
 215 networks (i.e., 10~50) for hopper tasks than PBRL and RORL that use 10  $Q$ -networks. (v) We  
 216 also include a basic SAC-10 algorithm as a baseline [2], which is an extension of SAC with 10  
 217  $Q$ -functions. To assign uniform adversarial attack budget on each dimension of observations, we  
 218 normalize the observations for SAC-10, EDAC and RORL. More details are provided in Appendix B

Table 1: Normalized average returns on Gym tasks, averaged over 4 random seeds. Part of the results are reported in the EDAC paper. **Top two scores for each task are highlighted.**

Task Name	BC	CQL	PBRL	SAC-10 (Reproduced)	EDAC (Paper)	EDAC-10 (Reproduced)	RORL (Ours)
halfcheetah-random	2.2±0.0	<b>31.3±3.5</b>	11.0±5.8	<b>29.0±1.5</b>	28.4±1.0	13.4 ± 1.1	28.6±0.6
halfcheetah-medium	43.2±0.6	46.9±0.4	57.9 ±1.5	64.9±1.3	<b>65.9±0.6</b>	64.1±1.1	<b>65.6±0.8</b>
halfcheetah-medium-expert	44.0±1.6	95.0±1.4	92.3±1.1	107.1±2.0	106.3±1.9	<b>107.2±1.0</b>	<b>107.8±0.7</b>
halfcheetah-medium-replay	37.6±2.1	45.3±0.3	45.1±8.0	<b>63.2±0.6</b>	61.3±1.9	60.1±0.3	<b>61.5±0.6</b>
halfcheetah-expert	91.8±1.5	97.3±1.1	92.4±1.7	104.9±0.9	<b>106.8±3.4</b>	104.0±0.8	<b>105.2±0.7</b>
hopper-random	3.7±0.6	5.3±0.6	<b>26.8±9.3</b>	25.9±9.6	25.3±10.4	16.9±10.1	<b>27.2±7.5</b>
hopper-medium	54.1±3.8	61.9±6.4	75.3±31.2	0.8±0.2	101.6±0.6	<b>103.6±0.2</b>	<b>105.0±0.3</b>
hopper-medium-expert	53.9±4.7	96.9±15.1	<b>110.8±0.8</b>	6.1±7.7	110.7±0.1	58.1±22.3	<b>112.7±0.2</b>
hopper-medium-replay	16.6±4.8	86.3±7.3	100.6±1.0	<b>102.9±0.9</b>	101.0±0.5	<b>102.8±0.3</b>	<b>102.1±0.4</b>
hopper-expert	107.7±9.7	106.5±9.1	<b>110.5±0.4</b>	1.1±0.5	110.1±0.1	77.0±43.9	<b>112.8±0.2</b>
walker2d-random	1.3±0.1	5.4±1.7	8.1±4.4	1.5±1.1	<b>16.6±7.0</b>	6.7±8.8	<b>21.4±0.2</b>
walker2d-medium	70.9±11.0	79.5±3.2	89.6±0.7	46.7±45.3	<b>92.5±0.8</b>	87.6±11.0	<b>102.4±1.4</b>
walker2d-medium-expert	90.1±13.2	109.1±0.2	110.1±0.3	<b>116.7±1.9</b>	114.7±0.9	115.4±0.5	<b>121.2±1.5</b>
walker2d-medium-replay	20.3±9.8	76.8±10.0	77.7±14.5	89.6±3.1	87.1±2.3	<b>94.0±1.2</b>	<b>90.4 ± 0.5</b>
walker2d-expert	108.7±0.2	109.3±0.1	108.3±0.3	1.2±0.7	<b>115.1±1.9</b>	57.8±55.7	<b>115.4 ± 0.5</b>
Average	49.7	70.2	74.4	50.8	<b>82.9</b>	71.2	<b>85.3</b>
Total	746.1	1052.8	1116.5	761.6	<b>1243.4</b>	1068.7	<b>1279.3</b>

## 219 6.1 Benchmark Results

220 We evaluate each method on Gym domain that includes three environments (Half Cheetah, Hopper,  
 221 and Walker2d) with four non-expert datasets (random, medium, medium-replay, and medium-expert).  
 222 The medium-replay dataset contains experiences collected in training a medium-level policy. The  
 223 random or medium dataset is generated by a single random or medium policy. The medium-  
 224 expert dataset is a mixture of medium and expert datasets. We set small perturbation scale  $\epsilon$  to  
 225 0.001/0.005/0.01 for halfcheetah/hopper/walker2d tasks respectively when training RORL and do not  
 226 include observation perturbation in the testing time.

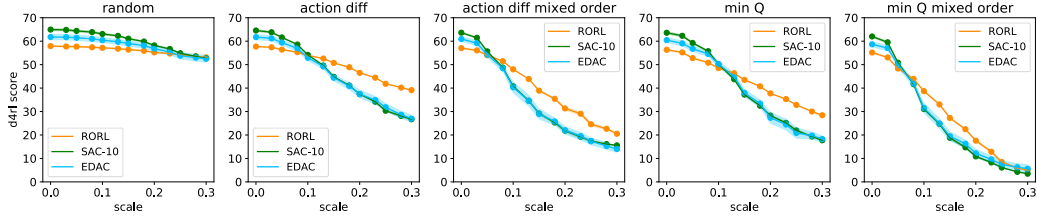
227 Table 1 reports the performance of the average normalized score with standard deviation. (i) SAC-10  
 228 is unstable in several Walker2d and Hopper tasks since the ensemble number is relatively small to  
 229 provide reliable uncertainties. (ii) EDAC solves this problem by using diversity constraints  
 230 while still requiring 10~50  $Q$ -networks to obtain reasonable performance. In contrast, RORL only  
 231 uses 10 ensemble  $Q$ -networks to outperform EDAC with 10~50 networks. (iii) PBRL chooses an  
 232 alternative OOD-sampling technique to reduce the ensemble numbers. According to the result, RORL  
 233 significantly outperforms PBRL with the same ensemble number. The reason is RORL additionally  
 234 uses conservative smoothing loss for both OOD states and actions, which improves the generalization  
 235 ability of the policy and  $Q$ -functions. We remark that RORL significantly improves over the current  
 236 SOTA results on walker2d and hopper tasks, which might be because the two tasks require a more  
 237 precise balancing of conservatism and generalization for better performance.

## 238 6.2 Adversarial Attack

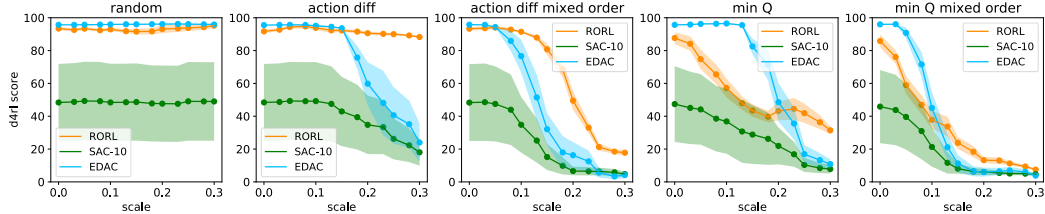
239 We adopt three attack methods, namely *random*, *action diff*, and *min Q* following prior works [64, 39].  
 240 Given perturbation scale  $\epsilon$ , the later two methods perform adversarial perturbation on observations  
 241 and are given access to the agent’s policy and value functions.

- 242 • *random* uniformly samples perturbed states in an  $l_\infty$  ball of norm  $\epsilon$ .
- 243 • *action diff* is an effective attack based on the agent’s policy and is proved to max-  
 244 imize an upper bound on the performance difference between perturbed and unperturbed  
 245 environments [64]. It directly finds states in an  $l_\infty$  ball of norm  $\epsilon$  to minimize:  
 246  $\min_{\hat{s} \in \mathbb{B}_d(s, \epsilon)} -D_J(\pi_\theta(\cdot|s) \parallel \pi_\theta(\cdot|\hat{s}))$ .
- 247 • *min Q* requires both the agent’s policy and value function to perform a relatively stronger  
 248 attack. The attacker finds a perturbed state to minimize the expected return of taking an  
 249 action from that state:  $\min_{\hat{s} \in \mathbb{B}_d(s, \epsilon)} Q(s, \pi_\theta(\hat{s}))$ . For ensemble-based algorithms,  $Q$  is set  
 250 as the mean of ensemble  $Q$  functions.

251 In our experiments, the two objectives of *action diff* and *min Q* are optimized via two ways. Specifi-  
 252 cally, we optimize the objectives through:



(a) Performance under attack on halfcheetah-medium-v2 dataset



(b) Performance under attack on walker2d-medium-v2 dataset

Figure 4: Figures (a) and (b) illustrate the performance of RORL, SAC-10, and EDAC under different types of attack and attack scales range  $[0, 0.3]$ . The curves are averaged over 3 seeds and smoothed with a window size of 3. The shaded region represents half a standard deviation.

- 253 (1) selecting the best state from uniformly sampled 50 states, which has the advantage of  
 254 simplicity and little computation cost.
- 255 (2) uniformly sampling 20 initial states and performing gradient descent for 10 steps with a step  
 256 size of  $\frac{1}{10}\epsilon$  from each initial state to find the best perturbed state. Note that we need to clip  
 257 the perturbed states within the  $l_\infty$  ball at the end of each optimization step. We name this  
 258 optimization method “mixed order” optimization.

259 We compare RORL with ensemble-based baselines EDAC and SAC-10 on halfcheetah-medium-v2  
 260 and walker2d-medium-v2 datasets. To handle large adversarial noise, we set the perturbation scale of  
 261 the policy and Q functions as  $[0.1, 0.05]$  for halfcheetah-medium-v2 and  $[0.1, 0.03]$  for walker2d-  
 262 medium-v2 in RORL’s training phase. More detailed description can be found in Appendix B. The  
 263 results are shown in Figure 4. In the results, RORL exhibits improved robustness than other baselines  
 264 under adversarial attacks and decreases the slowest as  $\epsilon$  increases. On the other hand, we find that  
 265 random attack is not effective for ensemble-based offline RL algorithms, and the “mixed order” attack  
 266 brings more significant performance drop than vanilla zero-order optimization.

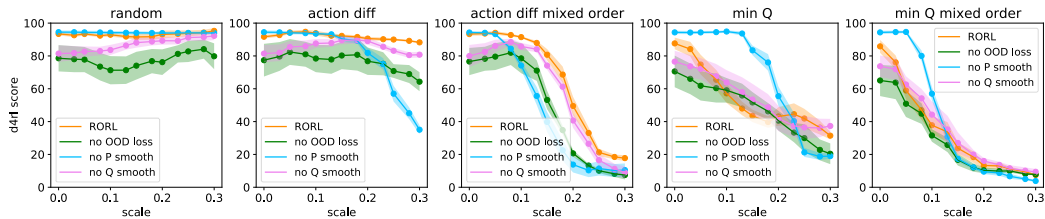


Figure 5: Ablation studies on walker2d-medium-v2 with varying perturbation scale. The curve is averaged across 3 random seeds and smoothed with a window size of 3. The shaded region represents half a standard deviation.

### 267 6.3 Ablations

268 We conducted ablation studies on walker2d-medium-v2 dataset to evaluate the importance of the  
 269 policy smoothing term, the Q function smoothing term and the penalty for OOD states in Figure  
 270 5. From the results, we can conclude that each loss contributes to the performance of RORL under  
 271 adversarial observation attacks. The policy smoothing loss is essential when the perturbation scale is  
 272 large, while the performance of “no P smooth” is higher than RORL when the scale is smaller than 0.1.  
 273 Besides, learning without OOD loss performs worse than RORL at almost all the perturbation scales.

274 In addition, Q smooth loss has the minimal impact, which is reasonable since the basic algorithm  
275 SAC-10 is based on 10 ensemble Q networks. Our smoothing technique may have a more significant  
276 impact on baselines with fewer Q networks, such as CQL in Section 3.

## 277 7 Related Works

278 **Offline RL** Research related to offline RL has experienced explosive growth in recent years. In  
279 model-free domain, offline RL methods focus on correcting the extrapolation error [13] in the off-  
280 policy algorithms. The natural idea is to regularize the learned policy near the dataset distribution [52,  
281 56, 33, 54, 60, 12]. For example, MARVIL reweights the policy with exponential advantage, which  
282 implicitly guarantees the policy within the KL-divergence neighborhood of the behavior policy.  
283 Another stream of model-free methods prevents the selection of OOD actions by penalizing their  
284 Q-value [26, 27, 2, 8]. With the ensemble Q networks and the additional loss term to diversify their  
285 gradients, EDAC [2] achieves SOTA performance in the D4RL benchmark. Instead of diversifying  
286 gradients, PBRL [4] proposes explicit training of the OOD actions, which achieves comparable  
287 performance with fewer ensemble networks. Inspired by EDAC and PBRL, we build our work upon  
288 ensemble networks, focusing more on smoothness over the state space.

289 Besides the surprising empirical results, theoretical analysis of offline reinforcement learning algo-  
290 rithms is of increasing interest [7, 22, 41, 58, 61]. Though the assumptions for the dataset vary in the  
291 different papers, they all suggest that pessimism and conservatism are necessary for offline RL. Our  
292 paper can be viewed as a robust extension to the previous work [22] (see Section 5 for details).

293 **Robust RL** The research line of robust RL can be traced back to  $H_\infty$ -control theory [57, 5],  
294 where policies are optimized to be well-performed in the worst possible deterministic environment.  
295 Depending on the definition, there are different streams of research on robust RL. As the extension  
296 of robust control to MDPs, Robust MDPs (RMDPs) [34, 19, 42, 17] are proposed to formulate the  
297 perturbation of transition probabilities for MDPs. Though some recent analyses with theoretical  
298 guarantees come out under specific assumptions for RMDP [66, 59, 29], there is currently no practical  
299 algorithm to solve RMDP in a large-scale problem, expect some linear approximation attempt [47].  
300 In online RL, domain randomization [49, 31] assumes the model uncertainty can be predefined in data  
301 collection by changing the setup of a simulator. However, it is not practical for offline RL. Robust  
302 Adversarial Reinforcement Learning (RARL) [39] and Noisy Robust Markov Decision Process  
303 (NR-MDP) [23] study the robust RL with the perturbed actions, showing that the policy robustness to  
304 adversarial or noisy actions can also induce robustness for model parameter changes.

305 The most related work to ours is SR<sup>2</sup>L [44], which shows policy smoothing can lead to significant  
306 performance improvements in the online setting. In contrast, we focus on the offline setting and  
307 tackle the potential overestimation of perturbed states. Another related work is S4RL [46], where the  
308 authors study different data augmentation methods to smooth observations in offline RL. Their result  
309 supports the necessity of state smoothing. More related works are given in Appendix D.

## 310 8 Conclusion

311 We propose Robust Offline Reinforcement Learning (RORL) to trade-off conservatism and robustness  
312 for offline RL. To achieve that, we introduce the conservative smoothing for the perturbed states  
313 while actively underestimating their values based on pessimistic bootstrapping to keep conservative.  
314 We show that RORL can achieve comparable performance with fewer ensemble Q networks than the  
315 previous methods in the offline RL benchmark. The result supports that our proposed conservative  
316 smoothing techniques benefit general offline RL by improving the generalization ability. In addition,  
317 we demonstrate that RORL is also robust to adversarial perturbations across the different types of  
318 attacks. We hope our work can promote the application of offline RL under real-world engineering  
319 conditions.

320 The main limitation of our method is that the adversarial state sampling slows down the computing  
321 process, which may be improved in future work.

## References

- 322
- 323 [1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear  
324 stochastic bandits. In *Advances in neural information processing systems*, volume 24, pages  
325 2312–2320, 2011.
- 326 [2] Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline  
327 reinforcement learning with diversified q-ensemble. *Advances in Neural Information Processing*  
328 *Systems*, 34, 2021.
- 329 [3] Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for  
330 reinforcement learning. In *International Conference on Machine Learning*, 2017.
- 331 [4] Chenjia Bai, Lingxiao Wang, Zhuoran Yang, Zhi-Hong Deng, Animesh Garg, Peng Liu, and  
332 Zhaoran Wang. Pessimistic bootstrapping for uncertainty-driven offline reinforcement learning.  
333 In *International Conference on Learning Representations*, 2022.
- 334 [5] Tamer Başar and Pierre Bernhard. *H-infinity optimal control and related minimax design*  
335 *problems: a dynamic game approach*. Springer Science & Business Media, 2008.
- 336 [6] Vahid Behzadan and Arslan Munir. Vulnerability of deep reinforcement learning to policy  
337 induction attacks. In *International Conference on Machine Learning and Data Mining in*  
338 *Pattern Recognition*, pages 262–275. Springer, 2017.
- 339 [7] Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement  
340 learning. In *International Conference on Machine Learning*, pages 1042–1051. PMLR, 2019.
- 341 [8] Ching-An Cheng, Tengyang Xie, Nan Jiang, and Alekh Agarwal. Adversarially trained actor  
342 critic for offline reinforcement learning. *arXiv preprint arXiv:2202.02446*, 2022.
- 343 [9] Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O Stanley, and Jeff Clune. First return,  
344 then explore. *Nature*, 590(7847):580–586, 2021.
- 345 [10] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for  
346 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 347 [11] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for  
348 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 349 [12] Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning.  
350 *Advances in Neural Information Processing Systems*, 34, 2021.
- 351 [13] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning  
352 without exploration. In *ICML*, 2019.
- 353 [14] Adam Gleave, Michael Dennis, Cody Wild, Neel Kant, Sergey Levine, and Stuart Russell.  
354 Adversarial policies: Attacking deep reinforcement learning. In *International Conference on*  
355 *Learning Representations*, 2019.
- 356 [15] Florin Gogianu, Tudor Berariu, Mihaela C Rosca, Claudia Clopath, Lucian Busoniu, and Razvan  
357 Pascanu. Spectral normalisation for deep reinforcement learning: an optimisation perspective.  
358 In *International Conference on Machine Learning*, pages 3734–3744. PMLR, 2021.
- 359 [16] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversar-  
360 ial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- 361 [17] Chin Pang Ho, Marek Petrik, and Wolfram Wiesemann. Fast Bellman Updates for Robust  
362 MDPs. In *Proceedings of the 35th International Conference on Machine Learning*, pages  
363 1979–1988. PMLR, 2018.
- 364 [18] Sandy Huang, Nicolas Papernot, Ian Goodfellow, Yan Duan, and Pieter Abbeel. Adversarial  
365 attacks on neural network policies. *arXiv preprint arXiv:1702.02284*, 2017.
- 366 [19] Garud N Iyengar. Robust dynamic programming. *Mathematics of Operations Research*,  
367 30(2):257–280, 2005.

- 368 [20] Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big  
369 sequence modeling problem. *Advances in neural information processing systems*, 34, 2021.
- 370 [21] Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement  
371 learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–  
372 2143. PMLR, 2020.
- 373 [22] Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In  
374 *International Conference on Machine Learning*, pages 5084–5096. PMLR, 2021.
- 375 [23] Parameswaran Kamalaruban, Yu-Ting Huang, Ya-Ping Hsieh, Paul Rolland, Cheng Shi, and  
376 Volkan Cevher. Robust reinforcement learning via adversarial training with langevin dynamics.  
377 *Advances in Neural Information Processing Systems*, 33:8127–8138, 2020.
- 378 [24] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel:  
379 Model-based offline reinforcement learning. In *NeurIPS*, 2020.
- 380 [25] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit  
381 q-learning. In *International Conference on Learning Representations*, 2021.
- 382 [26] Aviral Kumar, Justin Fu, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning  
383 via bootstrapping error reduction. In *NeurIPS*, 2019.
- 384 [27] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for  
385 offline reinforcement learning. In *NeurIPS*, 2020.
- 386 [28] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning:  
387 Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 388 [29] Jialian Li, Tongzheng Ren, Dong Yan, Hang Su, and Jun Zhu. Policy learning for robust markov  
389 decision process with a mismatched generative model. *arXiv preprint arXiv:2203.06587*, 2022.
- 390 [30] Yuzhe Ma, Xuezhou Zhang, Wen Sun, and Jerry Zhu. Policy poisoning in batch reinforcement  
391 learning and control. *Advances in Neural Information Processing Systems*, 32, 2019.
- 392 [31] Bhairav Mehta, Manfred Diaz, Florian Golemo, Christopher J Pal, and Liam Paull. Active  
393 domain randomization. In *Conference on Robot Learning*, pages 1162–1176. PMLR, 2020.
- 394 [32] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G  
395 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.  
396 Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 397 [33] Ashvin Nair, Murtaza Dalal, Abhishek Gupta, and Sergey Levine. Accelerating online rein-  
398 forcement learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.
- 399 [34] Arnab Nilim and Laurent Ghaoui. Robustness in markov decision problems with uncertain  
400 transition matrices. *Advances in neural information processing systems*, 16, 2003.
- 401 [35] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via  
402 bootstrapped dqn. In *NeurIPS*, 2016.
- 403 [36] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Anan-  
404 thram Swami. Practical black-box attacks against deep learning systems using adversarial  
405 examples. *arXiv preprint arXiv:1602.02697*, 1(2):3, 2016.
- 406 [37] Anay Pattanaik, Zhenyi Tang, Shuijing Liu, Gautham Bommannan, and Girish Chowdhary.  
407 Robust deep reinforcement learning with adversarial attacks. *arXiv preprint arXiv:1712.03632*,  
408 2017.
- 409 [38] Anay Pattanaik, Zhenyi Tang, Shuijing Liu, Gautham Bommannan, and Girish Chowdhary.  
410 Robust deep reinforcement learning with adversarial attacks. In *Proceedings of the 17th*  
411 *International Conference on Autonomous Agents and MultiAgent Systems*, pages 2040–2042,  
412 2018.

- 413 [39] Lerrel Pinto, James Davidson, Rahul Sukthankar, and Abhinav Gupta. Robust adversarial  
414 reinforcement learning. In *International Conference on Machine Learning*, pages 2817–2826.  
415 PMLR, 2017.
- 416 [40] Aravind Rajeswaran, Sarvjeet Ghotra, Balaraman Ravindran, and Sergey Levine. Epopt: Learn-  
417 ing robust neural network policies using model ensembles. *arXiv preprint arXiv:1610.01283*,  
418 2016.
- 419 [41] Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline  
420 reinforcement learning and imitation learning: A tale of pessimism. *Advances in Neural  
421 Information Processing Systems*, 34, 2021.
- 422 [42] Aurko Roy, Huan Xu, and Sebastian Pokutta. Reinforcement learning under model mismatch.  
423 *Advances in neural information processing systems*, 30, 2017.
- 424 [43] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Si-  
425 mon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering  
426 atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- 427 [44] Qianli Shen, Yan Li, Haoming Jiang, Zhaoran Wang, and Tuo Zhao. Deep reinforcement  
428 learning with robust and smooth policy. In *International Conference on Machine Learning*,  
429 pages 8707–8718. PMLR, 2020.
- 430 [45] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driess-  
431 che, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mas-  
432 tering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489,  
433 2016.
- 434 [46] Samarth Sinha, Ajay Mandlekar, and Animesh Garg. S4rl: Surprisingly simple self-supervision  
435 for offline reinforcement learning in robotics. In *Conference on Robot Learning*, pages 907–917.  
436 PMLR, 2022.
- 437 [47] Aviv Tamar, Huan Xu, and Shie Mannor. Scaling up robust mdps by reinforcement learning.  
438 *arXiv preprint arXiv:1306.6189*, 2013.
- 439 [48] Michael E Tipping and Christopher M Bishop. Probabilistic principal component analysis.  
440 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 61(3):611–622,  
441 1999.
- 442 [49] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel.  
443 Domain randomization for transferring deep neural networks from simulation to the real world.  
444 In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages  
445 23–30. IEEE, 2017.
- 446 [50] Eugene Vinitzky, Yuqing Du, Kanaad Parvate, Kathy Jang, Pieter Abbeel, and Alexandre Bayen.  
447 Robust reinforcement learning using adversarial populations. *arXiv preprint arXiv:2008.01825*,  
448 2020.
- 449 [51] Jianhao Wang, Wenzhe Li, Haozhe Jiang, Guangxiang Zhu, Siyuan Li, and Chongjie Zhang.  
450 Offline reinforcement learning with reverse model-based imagination. *Advances in Neural  
451 Information Processing Systems*, 34, 2021.
- 452 [52] Qing Wang, Jiechao Xiong, Lei Han, Han Liu, Tong Zhang, et al. Exponentially weighted  
453 imitation learning for batched historical data. *Advances in Neural Information Processing  
454 Systems*, 31, 2018.
- 455 [53] Ruosong Wang, Simon S Du, Lin F Yang, and Ruslan Salakhutdinov. On reward-free rein-  
456 forcement learning with linear function approximation. In *Advances in neural information  
457 processing systems*, 2020.
- 458 [54] Ziyu Wang, Alexander Novikov, Konrad Zolna, Josh S Merel, Jost Tobias Springenberg, Scott E  
459 Reed, Bobak Shahriari, Noah Siegel, Caglar Gulcehre, Nicolas Heess, et al. Critic regularized  
460 regression. *Advances in Neural Information Processing Systems*, 33:7768–7778, 2020.

- 461 [55] Fan Wu, Linyi Li, Chejian Xu, Huan Zhang, Bhavya Kailkhura, Krishnaram Kenthapadi, Ding  
462 Zhao, and Bo Li. Copa: Certifying robust policies for offline reinforcement learning against  
463 poisoning attacks. *arXiv preprint arXiv:2203.08398*, 2022.
- 464 [56] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement  
465 learning. *arXiv preprint arXiv:1911.11361*, 2019.
- 466 [57] Lihua Xie and Carlos E de Souza. Robust h/sub infinity/control for linear systems with norm-  
467 bounded time-varying uncertainty. In *29th IEEE Conference on Decision and Control*, pages  
468 1034–1035. IEEE, 1990.
- 469 [58] Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-  
470 consistent pessimism for offline reinforcement learning. *Advances in neural information  
471 processing systems*, 34, 2021.
- 472 [59] Wenhao Yang, Liangyu Zhang, and Zhihua Zhang. Towards theoretical understandings of  
473 robust markov decision processes: Sample complexity and asymptotics. *arXiv preprint  
474 arXiv:2105.03863*, 2021.
- 475 [60] Yiqin Yang, Xiaoteng Ma, Li Chenghao, Zewu Zheng, Qiyuan Zhang, Gao Huang, Jun Yang,  
476 and Qianchuan Zhao. Believe what you see: Implicit constraint approach for offline multi-agent  
477 reinforcement learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- 478 [61] Ming Yin, Yaqi Duan, Mengdi Wang, and Yu-Xiang Wang. Near-optimal offline reinforcement  
479 learning with linear representation: Leveraging variance information with pessimism. *arXiv  
480 preprint arXiv:2203.05804*, 2022.
- 481 [62] Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea  
482 Finn. Combo: Conservative offline model-based policy optimization. *Advances in Neural  
483 Information Processing Systems*, 34, 2021.
- 484 [63] Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Zou, Sergey Levine, Chelsea  
485 Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. In *NeurIPS*, 2020.
- 486 [64] Huan Zhang, Hongge Chen, Chaowei Xiao, Bo Li, Mingyan Liu, Duane Boning, and Cho-  
487 Jui Hsieh. Robust deep reinforcement learning against adversarial perturbations on state  
488 observations. *Advances in Neural Information Processing Systems*, 33:21024–21037, 2020.
- 489 [65] Xuezhou Zhang, Yiding Chen, Xiaojin Zhu, and Wen Sun. Robust policy gradient against  
490 strong data corruption. In *International Conference on Machine Learning*, pages 12391–12401.  
491 PMLR, 2021.
- 492 [66] Zhengqing Zhou, Zhengyuan Zhou, Qinxun Bai, Linhai Qiu, Jose Blanchet, and Peter Glynn.  
493 Finite-sample regret bound for distributionally robust offline tabular reinforcement learning. In  
494 *International Conference on Artificial Intelligence and Statistics*, pages 3331–3339. PMLR,  
495 2021.

## 496 Checklist

497 The checklist follows the references. Please read the checklist guidelines carefully for information on  
498 how to answer these questions. For each question, change the default [TODO] to [Yes], [No], or  
499 [N/A]. You are strongly encouraged to include a **justification to your answer**, either by referencing  
500 the appropriate section of your paper or providing a brief inline description. For example:

- 501 • Did you include the license to the code and datasets? [Yes] See Section ??.
- 502 • Did you include the license to the code and datasets? [No] The code and the data are  
503 proprietary.
- 504 • Did you include the license to the code and datasets? [N/A]

505 Please do not modify the questions and only use the provided macros for your answers. Note that the  
506 Checklist section does not count towards the page limit. In your paper, please delete this instructions  
507 block and only keep the Checklist section heading above along with the questions/answers below.

508 1. For all authors...

- 509 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
510 contributions and scope? [Yes]
- 511 (b) Did you describe the limitations of your work? [Yes]
- 512 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 513 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
514 them? [Yes]

515 2. If you are including theoretical results...

- 516 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 517 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix [A](#).

518 3. If you ran experiments...

- 519 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
520 mental results (either in the supplemental material or as a URL)? [Yes] See Sec [1](#) and  
521 Appendix [B](#).
- 522 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
523 were chosen)? [Yes] See Appendix [B](#).
- 524 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
525 ments multiple times)? [Yes] See Sec [6](#).
- 526 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
527 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix [B](#).

528 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 529 (a) If your work uses existing assets, did you cite the creators? [Yes] We cited D4RL [\[1\]](#)  
530 and EDAC [\[2\]](#) for their datasets and code.
- 531 (b) Did you mention the license of the assets? [Yes]
- 532 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
533 We included our code in the anonymized link.
- 534 (d) Did you discuss whether and how consent was obtained from people whose data you're  
535 using/curating? [Yes] Opensource code and dataset.
- 536 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
537 information or offensive content? [N/A]

538 5. If you used crowdsourcing or conducted research with human subjects...

- 539 (a) Did you include the full text of instructions given to participants and screenshots, if  
540 applicable? [N/A]
- 541 (b) Did you describe any potential participant risks, with links to Institutional Review  
542 Board (IRB) approvals, if applicable? [N/A]
- 543 (c) Did you include the estimated hourly wage paid to participants and the total amount  
544 spent on participant compensation? [N/A]