

# 1 Rebuttal Supplementary

## 1.1 Impact of Various $\alpha$ s

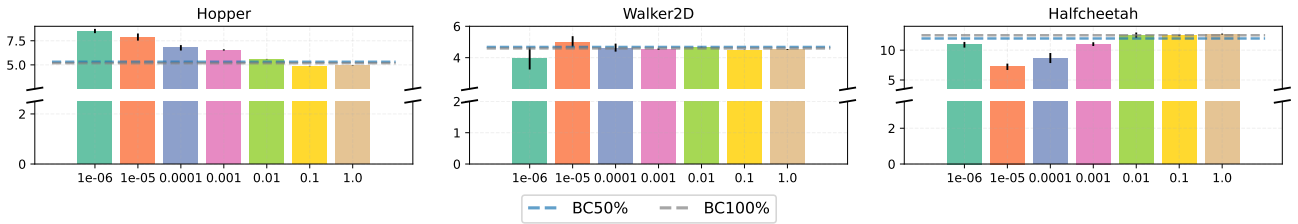


Figure 1: Comparison of different levels of the hyperparameter ( $\alpha$ ) ofROIDICE in locomotion environments using expert data quality. We average the scores and get  $\pm 2\times$  standard error with 5 seeds across 10 episodes.

## 1.2 Comparison of offline constrained RL algorithms

| Data quality | medium           | medium-expert    | expert           |
|--------------|------------------|------------------|------------------|
| ROIDICE      | 9.21 $\pm$ 0.49  | 8.29 $\pm$ 0.18  | 8.51 $\pm$ 0.25  |
| VOCE 50th    | 1.807 $\pm$ 0.63 | 1.33 $\pm$ 1.28  | 1.34 $\pm$ 0.67  |
| VOCE 80th    | 1.799 $\pm$ 0.62 | 1.34 $\pm$ 1.3   | 1.47 $\pm$ 0.76  |
| CPQ 50th     | 1.13 $\pm$ 0.58  | -0.06 $\pm$ 0.41 | -0.11 $\pm$ 0.71 |
| CPQ 80th     | -0.10 $\pm$ 0.64 | -0.15 $\pm$ 0.48 | -0.70 $\pm$ 0.19 |

Table 1: ROI ofROIDICE compared with offline constrained RL algorithms. We average the each scores and get  $\pm 2\times$  standard error with 5 seeds across 10 episodes.

## 1.3ROIDICE pseudocode

$$\mathcal{L}_{\nu_\phi} = \mathbb{E}_{s \sim p_0} [t(1 - \gamma)\nu(s)] + \mathbb{E}_{(s,a) \sim d_D} [w_{\nu,\mu,t}^*(s,a)(e_\nu(s,a) - \mu c(s,a)) - \alpha f(w_{\nu,\mu,t}^*(s,a), t)] \quad (1)$$

$$\mathcal{L}_\mu = \mathbb{E}_{(s,a) \sim d_D} [w_{\nu,\mu,t}^*(s,a)(e_\nu(s,a) - \mu c(s,a)) - \alpha f(w_{\nu,\mu,t}^*(s,a), t)] + \mu \quad (2)$$

$$\mathcal{L}_t = \mathbb{E}_{s \sim p_0} [t(1 - \gamma)\nu(s)] - \mathbb{E}_{(s,a) \sim d_D} [\alpha f(w_{\nu,\mu,t}^*(s,a), t)] \quad (3)$$

---

### Algorithm 1ROIDICE

---

**Input:** The offline dataset  $D$ , the initial state offline dataset  $p_0$ , parameterized Lagrangian multipliers  $\nu_\phi, \mu, t$ , and policy  $\pi_\theta$ .

**Output:** Optimal policy  $\pi_\theta^*$ .

Initialize all parameters.

**while** convergence **do**

    Update  $\nu_\phi, \mu, t$  with  $\mathcal{L}_{\nu_\phi}, \mathcal{L}_\mu, -\mathcal{L}_t$

    Update  $\theta$  with  $\mathcal{L}_\theta = -\mathbb{E}_{(s,a) \sim D} [\frac{1}{t} \cdot w_{\nu,\mu,t}^*(s,a) \cdot \log \pi_\theta(a|s)]$

**end while**

---

## 1.4ROI ofROIDICE in Safety RL environments

| Task           | CarGoal          | PointPush        |
|----------------|------------------|------------------|
| ROIDICE        | 0.160 $\pm$ 0.02 | 0.033 $\pm$ 0.01 |
| COptiDICE 50th | 0.119 $\pm$ 0.01 | 0.026 $\pm$ 0.01 |
| COptiDICE 80th | 0.115 $\pm$ 0.01 | 0.022 $\pm$ 0.01 |

Table 2: ROI ofROIDICE compared with offline constrained RL algorithms. We average the each scores and get  $\pm 2\times$  standard error with 5 seeds across 10 episodes.