

ROIL: Robust Offline Imitation Learning

Gersi Doko¹, Guang Yang², Daniel S. Brown² Marek Petrik¹

¹University of New Hampshire, ²University of Utah

Summary

Motivation

Learning from data in a robust offline way is important in many fields, like health care, robotics or finance.

Limitations of existing methods

Reliance on \hat{u}_e leads to covariate shift for off-policy datasets.

Inability to specify reliance on \hat{u}_e .

No guarantees of policy convergence to u_e .

Our contributions

New algorithm for robust offline imitation learning.

Guaranteed convergence to the optimal policy for tabular domains.

Flexibility to define the reliance on \hat{u}_e .

Markov Decision Process (MDP)

MDP : consists of a tuple $\langle \mathcal{S}, \mathcal{A}, p, r \rangle$

- State space $\mathcal{S} = \{ \text{mild}, \text{moderate}, \text{severe} \}$

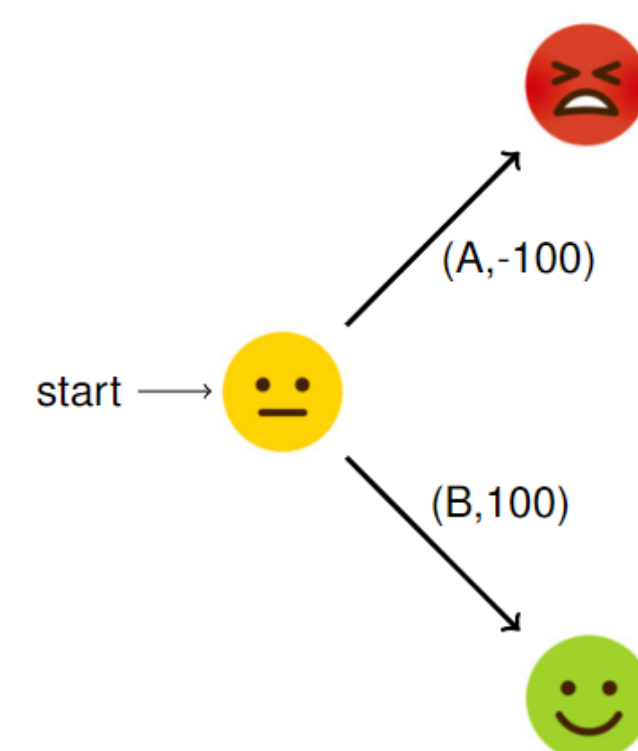


- Action space $\mathcal{A} = \{A, B\}$



- Transition probability p : $p(\text{moderate}, \text{pill}, \text{severe}) = 1$

- Reward function r : $r(\text{moderate}, \text{pill}) = 100$



Inverse Reinforcement Learning (IRL)

Given a dataset of expert demonstrations and an MDP model without a reward function, IRL aims to learn a policy π from a dataset of expert demonstrations $D_e = \{s_i, a_i, r_i\}_{i=0}^N$.

$$\rho(\pi, r) = \lim_{T \rightarrow \infty} \mathbb{E}^{\pi, p_0} \left[\sum_{t=0}^T \gamma^t r(\tilde{s}_t, \pi(\tilde{s}_t)) \right]$$

Optimal policy π^*

$$\pi^* = \arg \min_{\pi \in \Pi} \max_{\pi_e \in \Pi} \max_{r \in \mathcal{R}} \rho(\pi_e, r) - \rho(\pi, r)$$

Previous Work

LPAL
MILO
GAIL
BROIL

ROIL LP

$$\begin{aligned} \min_{t \in \mathbb{R}, u \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}} \quad & t \\ \text{s.t.} \quad & t \geq -u^T \Phi w + \max_{v \in \Upsilon} v^T \Phi w, \quad \forall w \in \text{ext}(\mathcal{W}), \\ & u \in \Upsilon, \end{aligned}$$

Gridworld Results

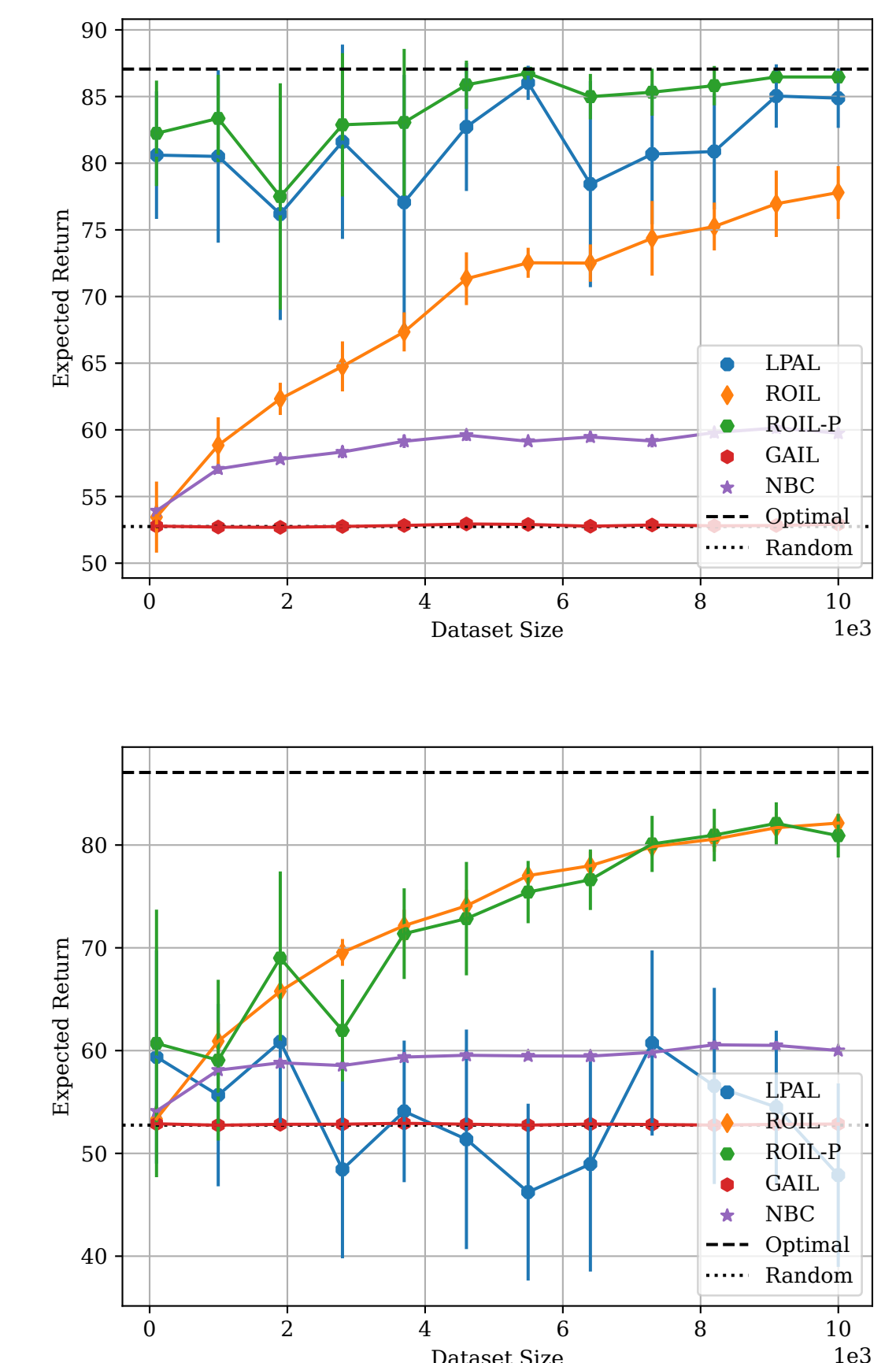


Figure: On-policy and off-policy returns (respectively) for 40x40 Gridworld.

Regret Results

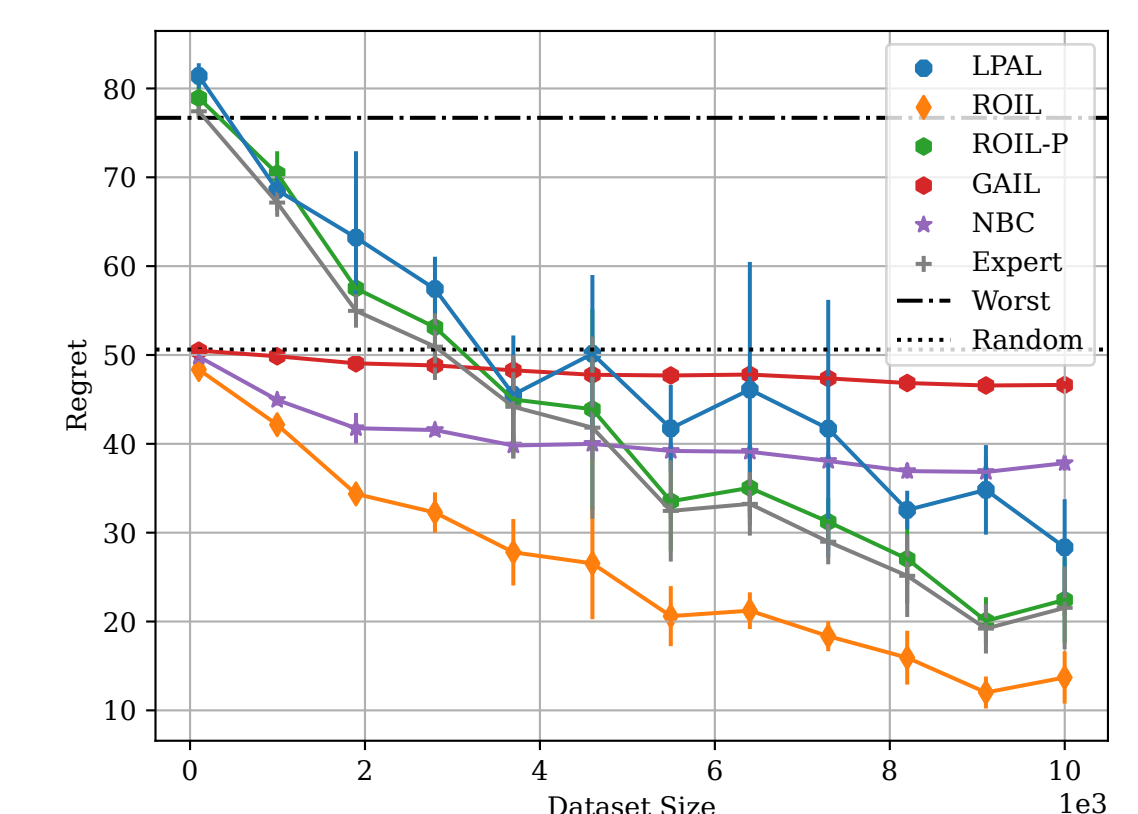


Figure: On-policy regret for 40x40 Gridworld.