

Adversarially Robust Imitation Learning

Anonymous Author(s)

Affiliation

Address

email

1 A Proof of Theorem 4.1

2 *Proof.* Since Alg. 1 runs no-regret online learner to update π on the sequence of loss functions
3 $\{\ell_t(\pi)\}$, we must have:

$$\sum_{t=0}^T \ell_t(\pi_t) - \min_{\pi \in \Pi} \sum_{t=0}^T \ell_t(\pi) \leq o(T). \quad (1)$$

4 Add and subtract $\sum_{t=0}^T \ell_t(\pi^e)$ on the left hand side of the above inequality, we have:

$$\sum_{t=0}^T \ell_t(\pi_t) - \sum_{t=0}^T \ell_t(\pi^e) \leq o(T) + \left(\min_{\pi \in \Pi} \sum_{t=0}^T \ell_t(\pi) - \sum_{t=0}^T \ell_t(\pi^e) \right) \quad (2)$$

5 Since we operate under the realizability setting, the term inside the parenthesis on the RHS of the
6 above inequality is guaranteed to be less than or equal to zero. Hence, the above inequality simplifies
7 to:

$$\sum_{t=0}^T \ell_t(\pi_t) \leq \sum_{t=0}^T \ell_t(\pi^e) + o(T). \quad (3)$$

8 For $\ell_t(\pi^e)$, using the definition of ℓ_t , we see:

$$\ell_t(\pi^e) = \mathbb{E}_{s \sim d_{\pi_t \circ f_t}} [\mathbb{E}_{a \sim \pi^e(f_t(s))} [\ell(a, \pi^e(s))]] = \mathbb{E}_{s \sim d_{\pi_t \circ f_t}} [\ell(\pi^e(f_t(s)), \pi^e(s))] = 0, \quad (4)$$

9 where we use the expert robustness definition above (i.e., $\pi^e(f(s)) = \pi^e(s)$ for all s and f). Hence,
10 we have:

$$\sum_{t=0}^T \ell_t(\pi_t) \leq T\epsilon_e + o(T). \quad (5)$$

11 Now we need to lower bound $\ell_t(\pi)$. Using the definition of ℓ_t again, we have:

$$\max_f \mathcal{L}(\pi_{i^*}, f) \leq \mathcal{L}(\pi_{i^*}, f_{i^*}) + \epsilon_{rl} \leq \frac{1}{T} \sum_t \mathcal{L}(\pi_t, f_t) + \epsilon_{rl} = \frac{1}{T} \sum_t \ell_t(\pi_t) + \epsilon_{rl}, \quad (6)$$

12 where the first inequality comes from the assumption that the RL solver returns a ϵ_{rl} near-optimal
13 solution.

14 Combine all the results above together, we get:

$$\max_f \mathcal{L}(\pi_{i^*}, f) \leq \epsilon_{rl} + \epsilon_e + o(T)/T. \quad (7)$$

15 Hence as T approaches to ∞ , we have that the long-term prediction loss of the learned policy π_{i^*}
16 under the worst possible adversarial attack from \mathcal{F} is upper bounded by ϵ_{rl} —the error introduced
17 from optimizing f_t .

18 Under agnostic setting, it is not guaranteed that there exists $a = \pi^e(s)$, where $a \sim \pi(s)$ and $\pi \in \Pi$
19 for $\forall s \in \mathcal{S}$. However, it can be assumed that the error between a and $\pi^e(s)$ is bounded by a small
20 number ϵ_a , namely $\mathbb{E}_{a \sim \pi(s)} [a \neq \pi^e(s)] \leq \epsilon_a$. Hence Equation 4 is modified as:

$$\ell_t(\pi^e) = \mathbb{E}_{s \sim d_{\pi_t \circ f_t}} [\mathbb{E}_{a \sim \pi^e(f_t(s))} [\ell(a, \pi^e(s))]] \leq \epsilon_a. \quad (8)$$

21 Plugging this to Equation 3, we have:

$$\sum_{t=0}^T \ell_t(\pi_t) \leq T(\epsilon_e + \epsilon_a) + o(T). \quad (9)$$

22 Combining the results with Equation 6, we get

$$\max_f \mathcal{L}(\pi_{i^*}, f) \leq \epsilon_{rl} + \epsilon_e + \epsilon_a + o(T)/T. \quad (10)$$

23 Similarly to the realizability setting, as T approaches to ∞ , the long-term prediction loss of the
 24 learned policy π_{i^*} under the worst possible adversarial attack from \mathcal{F} is still upper bounded.

25

□

26 B Experiments in Detail

27 B.0.1 Experiment Setting

28 Table 1 shows ϵ in each settings.

29 And as figure 1, the perturbation added to each state is bounded to a small value such that it's
 30 imperceptible but can lead to significant performance decrease for the student network before trained
 31 by ARIL.

Table 1: ϵ for Each Attack			
Environments	Ant-v2	HalfCheetah-v2	Swimmer-v2
sensory IL	0.5	0.04	0.25
physical	0.01	0.024	0.014

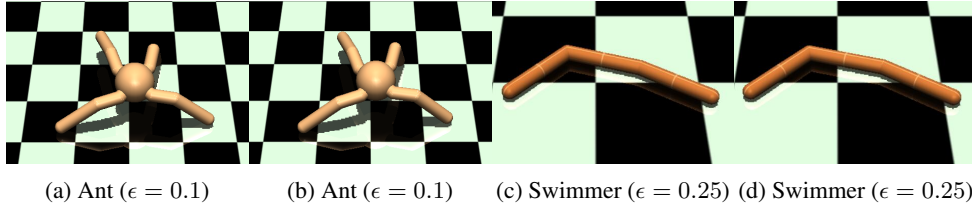


Figure 1: Robot observation before (a,c) and after (b,d) attack.

32 B.0.2 Training Details

33 Both attacker and student are trained using Adam optimizer at learning rate of 0.001.

34 At the start of ARIL algorithm, we initialized the buffer for DAgger algorithm with size $5M$ timesteps.
 35 In each attack stage, it takes $1e6$ timesteps for the attacker to collect the trajectories each time. Then
 36 at each defense stage, we collected 40 trajectories of student under attack and labeled them with
 37 expert actions into the DAgger buffer.

Table 2: Hyperparameters for training attacker

entropy loss coefficient	0.01
value loss coefficient	0.5
clip range	0.2