

## APPENDIX

In the appendix, we provide additional details of the theories and the experiments. The contents of this appendix are as follows.

- In Appendix A, we discuss some details of PAGAR and *MinimaxRegret* that were omitted in Section 4. We briefly introduce some necessary preliminaries in Appendix A.1. Then we derive a Theorem 5 to support Table 1 in Appendix A.2. The proofs for Theorem 1 and 1, Corollary 2 and 3, Theorem 1 are in Appendix A.4. The details of Example 1 are in Appendix A.5.
- In Appendix B, we provide some details of Imitation Learning with PAGAR that were omitted in Section 5. We prove Theorem 4 in Appendix B.1. Then we derive the objective functions in Appendix B.2. Some details of Algorithm 1 will be explained in Appendix B.3.
- In Appendix C, we provide some experimental details and additional results.

## A REWARD DESIGN WITH PAGAR

This paper does not aim to resolve the ambiguity problem in IRL but provides a way to circumvent it so that reward ambiguity does not lead to reward misalignment in IRL-based IL. PAGAR, the semi-supervised reward design paradigm proposed in this paper, tackles this problem from the perspective of semi-supervised reward design. But the nature of PAGAR is distinct from IRL and IL: assume that a set of reward functions is available for some underlying task, where some of those reward functions align with the task while others are misaligned, PAGAR provides a solution for selecting reward functions to train a policy that successfully performs the task, without knowing which reward function aligns with the task. Our research demonstrates that policy training with PAGAR is equivalent to learning a policy to maximize an affine combination of utilities measured under a distribution of the reward functions in the reward function set. With this understanding of PAGAR, we integrate it with IL to illustrate its advantages.

### A.1 SEMI-SUPERVISED REWARD DESIGN

Designing a reward function can be thought as deciding an ordering of policies. We adopt a concept, called *total domination*, from unsupervised environment design [Dennis et al. (2020)], and re-interpret this concept in the context of reward design. In this paper, we suppose that a function  $U_r(\pi)$  is given to measure the performance of a policy. While the measurement of policy performance can vary depending on the free variable  $r$ , *total dominance* can be viewed as an invariance regardless of such dependency.

**Definition 3** (Total Domination). A policy,  $\pi_1$ , is totally dominated by some policy  $\pi_2$  w.r.t a reward function set  $R$ , if for every pair of reward functions  $r_1, r_2 \in R$ ,  $U_{r_1}(\pi_1) < U_{r_2}(\pi_2)$ .

If  $\pi_1$  totally dominate  $\pi_2$  w.r.t  $R$ ,  $\pi_2$  can be regarded as being unconditionally better than  $\pi_1$ . In other words, the two sets  $\{U_r(\pi_1) | r \in R\}$  and  $\{U_r(\pi_2) | r \in R\}$  are disjoint, such that  $\sup\{U_r(\pi_1) | r \in R\} < \inf\{U_r(\pi_2) | r \in R\}$ . Conversely, if a policy  $\pi$  is not totally dominated by any other policy, it indicates that for any other policy, say  $\pi_2$ ,  $\sup\{U_r(\pi_1) | r \in R\} \geq \inf\{U_r(\pi_2) | r \in R\}$ .

**Definition 4.** A reward function set  $R$  aligns with an ordering  $\prec_R$  among policies such that  $\pi_1 \prec_R \pi_2$  if and only if  $\pi_1$  is totally dominated by  $\pi_2$  w.r.t.  $R$ .

Especially, designing a reward function  $r$  is to establish an ordering  $\prec_{\{r\}}$  among policies. Total domination can be extended to policy-conditioned reward design, where the reward function  $r$  is selected by following a decision rule  $\mathcal{R}(\pi)$  such that  $\sum_{r \in R} \mathcal{R}(\pi)(r) = 1$ . We let  $\mathcal{U}_{\mathcal{R}}(\pi) = \sum_{r \in R} \mathcal{R}(\pi)(r) \cdot U_r(\pi)$  be an affine combination of  $U_r(\pi)$ 's with its coefficients specified by  $\mathcal{R}(\pi)$ .

**Definition 5.** A policy conditioned decision rule  $\mathcal{R}$  is said to prefer a policy  $\pi_1$  to another policy  $\pi_2$ , which is notated as  $\pi_1 \prec^{\mathcal{R}} \pi_2$ , if and only if  $\mathcal{U}_{\mathcal{R}}(\pi_1) < \mathcal{U}_{\mathcal{R}}(\pi_2)$ .

Making a decision rule for selecting reward functions from a reward function set to respect the total dominance w.r.t this reward function set is an unsupervised learning problem, where no additional

external supervision is provided. If considering expert demonstrations as a form of supervision and using it to constrain the set  $R_E$  of reward function via IRL, the reward design becomes semi-supervised.

## A.2 SOLUTION TO THE MINIMAXREGRET

In Table [1](#), we mentioned that solving  $\text{MinimaxRegret}(R_E)$  is equivalent to finding an optimal policy  $\pi^*$  to maximize a  $\mathcal{U}_{R_E}(\pi)$  under a decision rule  $\mathcal{R}_E$ . Without loss of generality, we use  $R$  instead of  $R_E$  in our subsequent analysis, because solving  $\text{MinimaxRegret}(R)$  does not depend on whether there are constraints for  $R$ . In order to show such an equivalence, we follow the same routine as in [Dennis et al. \(2020\)](#), and start by introducing the concept of *weakly total domination*.

**Definition 6** (Weakly Total Domination). A policy  $\pi_1$  is *weakly totally dominated* w.r.t a reward function set  $R$  by some policy  $\pi_2$  if and only if for any pair of reward function  $r_1, r_2 \in R$ ,  $U_{r_1}(\pi_1) \leq U_{r_2}(\pi_2)$ .

Note that a policy  $\pi$  being totally dominated by any other policy is a sufficient but not necessary condition for  $\pi$  being weakly totally dominated by some other policy. A policy  $\pi_1$  being weakly totally dominated by a policy  $\pi_2$  implies that  $\sup\{U_r(\pi_1) | r \in R\} \leq \inf\{U_r(\pi_2) | r \in R\}$ . We assume that there does not exist a policy  $\pi$  that weakly totally dominates itself, which could happen if and only if  $U_r(\pi)$  is a constant. We formalize this assumption as the following.

**Assumption 1.** For the given reward set  $R$  and policy set  $\Pi$ , there does not exist a policy  $\pi$  such that for any two reward functions  $r_1, r_2 \in R$ ,  $U_{r_1}(\pi) = U_{r_2}(\pi)$ .

This assumption makes weak total domination a non-reflexive relation. It is obvious that weak total domination is transitive and asymmetric. Now we show that successive weak total domination will lead to total domination.

**Lemma 1.** *for any three policies  $\pi_1, \pi_2, \pi_3 \in \Pi$ , if  $\pi_1$  is weakly totally dominated by  $\pi_2$ ,  $\pi_2$  is weakly totally dominated by  $\pi_3$ , then  $\pi_3$  totally dominates  $\pi_1$ .*

*Proof.* According to the definition of weak total domination,  $\max_{r \in R} U_r(\pi_1) \leq \min_{r \in R} U_r(\pi_2)$  and  $\max_{r \in R} U_r(\pi_2) \leq \min_{r \in R} U_r(\pi_3)$ . If  $\pi_1$  is weakly totally dominated but not totally dominated by  $\pi_3$ , then  $\max_{r \in R} U_r(\pi_1) = \min_{r \in R} U_r(\pi_3)$  must be true. However, it implies  $\min_{r \in R} U_r(\pi_2) = \max_{r \in R} U_r(\pi_2)$ , which violates Assumption [1](#). We finish the proof.  $\square$

**Lemma 2.** *For the set  $\Pi_{\neg wtd} \subseteq \Pi$  of policies that are not weakly totally dominated by any other policy in the whole set of policies w.r.t a reward function set  $R$ , there exists a range  $\mathbb{U} \subseteq \mathbb{R}$  such that for any policy  $\pi \in \Pi_{\neg wtd}$ ,  $\mathbb{U} \subseteq [\min_{r \in R} U_r(\pi), \max_{r \in R} U_r(\pi)]$ .*

*Proof.* For any two policies  $\pi_1, \pi_2 \in \Pi_{\neg wtd}$ , it cannot be true that  $\max_{r \in R} U_r(\pi_1) = \min_{r \in R} U_r(\pi_2)$  nor  $\min_{r \in R} U_r(\pi_1) = \max_{r \in R} U_r(\pi_2)$ , because otherwise one of the policies weakly totally dominates the other. Without loss of generalization, we assume that  $\max_{r \in R} U_r(\pi_1) > \min_{r \in R} U_r(\pi_2)$ . In this case,  $\max_{r \in R} U_r(\pi_2) > \min_{r \in R} U_r(\pi_1)$  must also be true, otherwise  $\pi_1$  weakly totally dominates  $\pi_2$ . Inductively,  $\min_{\pi \in \Pi_{\neg wtd}} \max_{r \in R} U_r(\pi) > \max_{\pi \in \Pi_{\neg wtd}} \min_{r \in R} U_r(\pi)$ . Letting  $ub = \min_{\pi \in \Pi_{\neg wtd}} \max_{r \in R} U_r(\pi)$  and  $lb = \max_{\pi \in \Pi_{\neg wtd}} \min_{r \in R} U_r(\pi)$ , any  $\mathbb{U} \subseteq [lb, ub]$  shall support the assertion. We finish the proof.  $\square$

**Lemma 3.** *For a reward function set  $R$ , if a policy  $\pi \in \Pi$  is weakly totally dominated by some other policy in  $\Pi$  and there exists a subset  $\Pi_{\neg wtd} \subseteq \Pi$  of policies that are not weakly totally dominated by any other policy in  $\Pi$ , then  $\max_{r \in R} U_r(\pi) < \min_{\pi' \in \Pi_{\neg wtd}} \max_{r \in R} U_r(\pi')$*

*Proof.* If  $\pi_1$  is weakly totally dominated by a policy  $\pi_2 \in \Pi$ , then  $\min_{r \in R} U_r(\pi_2) = \max_{r \in R} U_r(\pi)$ . If  $\max_{r \in R} U_r(\pi) \geq \min_{\pi' \in \Pi_{\neg wtd}} \max_{r \in R} U_r(\pi')$ , then  $\min_{r \in R} U_r(\pi_2) \geq \min_{\pi' \in \Pi_{\neg wtd}} \max_{r \in R} U_r(\pi')$ , making at

least one of the policies in  $\Pi_{\neg wtd}$  being weakly totally dominated by  $\pi_2$ . Hence,  $\max_{r \in R} U_r(\pi) < \min_{\pi' \in \Pi_{\neg wtd}} \max_{r \in R} U_r(\pi')$  must be true.  $\square$

Given a policy  $\pi$  and a reward function  $r$ , the regret is represented as Eq[8]

$$Regret(\pi, r) := \max_{\pi'} U_r(\pi') - U_r(\pi) \quad (8)$$

Then we represent the *MinimaxRegret*( $R$ ) problem in Eq[9]

$$MinimaxRegret(R) := \arg \min_{\pi \in \Pi} \left\{ \max_{r \in R} Regret(\pi, r) \right\} \quad (9)$$

We denote as  $r_\pi^* \in R$  the reward function that maximizes  $U_r(\pi)$  among all the  $r$ 's that achieve the maximization in Eq[9]. Formally,

$$r_\pi^* \in \arg \max_{r \in R} U_r(\pi) \quad s.t. \ r \in \arg \max_{r' \in R} Regret(\pi, r') \quad (10)$$

Then *MinimaxRegret* can be defined as minimizing the worst-case regret as in Eq[9]. Next, we want to show that for some decision rule  $\mathcal{R}$ , the set of optimal policies which maximizes  $\mathcal{U}_{\mathcal{R}}$  are the solutions to *MinimaxRegret*( $R$ ). Formally,

$$MinimaxRegret(R) = \arg \max_{\pi \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi) \quad (11)$$

We design  $\mathcal{R}$  by letting  $\mathcal{R}(\pi) := \overline{\mathcal{R}}(\pi) \cdot \delta_{r_\pi^*} + (1 - \overline{\mathcal{R}}(\pi)) \cdot \underline{\mathcal{R}}(\pi)$  where  $\underline{\mathcal{R}} : \Pi \rightarrow \Delta(R)$  is a policy conditioned distribution over reward functions,  $\delta_{r_\pi^*}$  be a delta distribution centered at  $r_\pi^*$ , and  $\overline{\mathcal{R}}(\pi)$  is a coefficient. We show how to design  $\underline{\mathcal{R}}$  by using the following lemma.

**Lemma 4.** *Given that the reward function set is  $R$ , there exists a decision rule  $\underline{\mathcal{R}} : \Pi \rightarrow \Delta(R)$  which guarantees that: 1) for any policy  $\pi$  that is not weakly totally dominated by any other policy in  $\Pi$ , i.e.,  $\pi \in \Pi_{\neg wtd} \subseteq \Pi$ ,  $\mathcal{U}_{\underline{\mathcal{R}}}(\pi) \equiv c$  where  $c = \max_{\pi' \in \Pi_{\neg wtd}} \min_{r \in R} U_r(\pi')$ ; 2) for any  $\pi$  that is weakly totally dominated by some policy but not totally dominated by any policy,  $\mathcal{U}_{\underline{\mathcal{R}}}(\pi) = \max_{r \in R} U_r(\pi)$ ; 3) if  $\pi$  is totally dominated by some other policy,  $\overline{\mathcal{R}}(\pi)$  is a uniform distribution.*

*Proof.* Since the description of  $\underline{\mathcal{R}}$  for the policies in condition 2) and 3) are self-explanatory, we omit the discussion on them. For the none weakly totally dominated policies in condition 1), having a constant  $\mathcal{U}_{\underline{\mathcal{R}}}(\pi) \equiv c$  is possible if and only if for any policy  $\pi \in \Pi_{\neg wtd}$ ,  $c \in [\min_{r \in R} U_r(\pi'), \max_{r \in R} U_r(\pi')]$ .

As mentioned in the proof of Lemma 2,  $c$  can exist within  $[\min_{r \in R} U_r(\pi), \max_{r \in R} U_r(\pi)]$ . Hence,  $c = \max_{\pi' \in \Pi_{\neg wtd}} \min_{r \in R} U_r(\pi')$  is a valid assignment.  $\square$

Then by letting  $\overline{\mathcal{R}}(\pi) := \frac{Regret(\pi, r_\pi^*)}{c - U_{r_\pi^*}(\pi)}$ , we have the following theorem.

**Theorem 5.** *By letting  $\mathcal{R}(\pi) := \overline{\mathcal{R}}(\pi) \cdot \delta_{r_\pi^*} + (1 - \overline{\mathcal{R}}(\pi)) \cdot \underline{\mathcal{R}}(\pi)$  with  $\overline{\mathcal{R}}(\pi) := \frac{Regret(\pi, r_\pi^*)}{c - U_{r_\pi^*}(\pi)}$  and any  $\underline{\mathcal{R}}$  that satisfies Lemma 4*

$$MinimaxRegret(R) = \arg \max_{\pi \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi) \quad (12)$$

*Proof.* If a policy  $\pi \in \Pi$  is totally dominated by some other policy, since there exists another policy with larger  $\mathcal{U}_{\mathcal{R}}$ ,  $\pi$  cannot be a solution to  $\arg \max_{\pi \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi)$ . Hence, there is no need for further discussion on totally dominated policies. We discuss the none weakly totally dominated policies

and the weakly totally dominated but not totally dominated policies (shortened to "weakly totally dominated" from now on) respectively. First we expand  $\arg \max_{\pi \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi)$  as in Eq 13.

$$\begin{aligned}
& \arg \max_{\pi \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi) \\
&= \arg \max_{\pi \in \Pi} \sum_{r \in R} \mathcal{R}(\pi)(r) \cdot U_r(\pi) \\
&= \arg \max_{\pi \in \Pi} \frac{\text{Regret}(\pi, r_{\pi}^*) \cdot U_{r_{\pi}^*}(\pi) + (\mathcal{U}_{\mathcal{R}}(\pi) - U_{r_{\pi}^*}(\pi) - \text{Regret}(\pi, r_{\pi}^*)) \cdot \mathcal{U}_{\mathcal{R}}(\pi)}{c - U_{r_{\pi}^*}(\pi)} \\
&= \arg \max_{\pi \in \Pi} \frac{(\mathcal{U}_{\mathcal{R}}(\pi) - U_{r_{\pi}^*}(\pi)) \cdot \mathcal{U}_{\mathcal{R}}(\pi) - (\mathcal{U}_{\mathcal{R}}(\pi) - U_{r_{\pi}^*}(\pi)) \cdot \text{Regret}(\pi, r_{\pi}^*)}{c - U_{r_{\pi}^*}(\pi)} \\
&= \arg \max_{\pi \in \Pi} \frac{\mathcal{U}_{\mathcal{R}}(\pi) - U_{r_{\pi}^*}(\pi)}{c - U_{r_{\pi}^*}(\pi)} \cdot \mathcal{U}_{\mathcal{R}}(\pi) - \text{Regret}(\pi, r_{\pi}^*) \tag{13}
\end{aligned}$$

1) For the none weakly totally dominated policies, since by design  $\mathcal{U}_{\mathcal{R}} \equiv c$ , Eq 13 is equivalent to  $\arg \max_{\pi \in \Pi_1} -\text{Regret}(\pi, r_{\pi}^*)$  which exactly equals  $\text{MinimaxRegret}(R)$ . Hence, the equivalence holds among the none weakly totally dominated policies. Furthermore, if a none weakly totally dominated policy  $\pi \in \Pi_{\neg wtd}$  achieves optimality in  $\text{MinimaxRegret}(R)$ , its  $\mathcal{U}_{\mathcal{R}}(\pi)$  is also no less than any weakly totally dominated policy. Because according to Lemma 3, for any weakly totally dominated policy  $\pi_1$ , its  $\mathcal{U}_{\mathcal{R}}(\pi_1) \leq c$ , hence  $\frac{\mathcal{U}_{\mathcal{R}}(\pi) - U_{r_{\pi}^*}(\pi)}{c - U_{r_{\pi}^*}(\pi)} \cdot \mathcal{U}_{\mathcal{R}}(\pi) \leq c$ . Since  $\text{Regret}(\pi, r_{\pi}^*) \leq \text{Regret}(\pi_1, r_{\pi_1}^*)$ ,  $\mathcal{U}_{\mathcal{R}}(\pi) \geq \mathcal{U}_{\mathcal{R}}(\pi_1)$ . Therefore, we can assert that if a none weakly totally dominated policy  $\pi$  is a solution to  $\text{MinimaxRegret}(R)$ , it is also a solution to  $\arg \max_{\pi \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi)$ . Additionally, to prove that if a none weakly totally dominated policy  $\pi$  is a solution to  $\arg \max_{\pi' \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi')$ , it is also a solution to  $\text{MinimaxRegret}(R)$ , it is only necessary to prove that  $\pi$  achieve no larger regret than all the weakly totally dominated policies. But we delay the proof to 2).

2) If a policy  $\pi$  is weakly totally dominated and is a solution to  $\text{MinimaxRegret}(R)$ , we show that it is also a solution to  $\arg \max_{\pi \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi)$ , i.e., its  $\mathcal{U}_{\mathcal{R}}(\pi)$  is no less than that of any other policy.

We start by comparing with non weakly totally dominated policy. for any weakly totally dominated policy  $\pi_1 \in \text{MinimaxRegret}(R)$ , it must hold true that  $\text{Regret}(\pi_1, r_{\pi_1}^*) \leq \text{Regret}(\pi_2, r_{\pi_2}^*)$  for any  $\pi_2 \in \Pi$  that weakly totally dominates  $\pi_1$ . However, it also holds that  $\text{Regret}(\pi_2, r_{\pi_2}^*) \leq \text{Regret}(\pi_1, r_{\pi_1}^*)$  due to the weak total domination. Therefore,  $\text{Regret}(\pi_1, r_{\pi_1}^*) = \text{Regret}(\pi_2, r_{\pi_2}^*) = \text{Regret}(\pi_1, r_{\pi_2}^*)$ , implying that  $\pi_2$  is also a solution to  $\text{MinimaxRegret}(R)$ . It also implies that  $U_{r_{\pi_2}^*}(\pi_1) = U_{r_{\pi_2}^*}(\pi_2) \geq U_{r_{\pi_1}^*}(\pi_1)$  due to the weak total domination. However, by definition  $U_{r_{\pi_1}^*}(\pi_1) \geq U_{r_{\pi_2}^*}(\pi_1)$ . Hence,  $U_{r_{\pi_1}^*}(\pi_1) = U_{r_{\pi_2}^*}(\pi_1) = U_{r_{\pi_2}^*}(\pi_2)$  must hold. Now we discuss two possibilities: a) there exists another policy  $\pi_3$  that weakly totally dominates  $\pi_2$ ; b) there does not exist any other policy that weakly totally dominates  $\pi_2$ . First, condition a) cannot hold. Because inductively it can be derived  $U_{r_{\pi_1}^*}(\pi_1) = U_{r_{\pi_2}^*}(\pi_1) = U_{r_{\pi_2}^*}(\pi_2) = U_{r_{\pi_3}^*}(\pi_3)$ , while Lemma 1 indicates that  $\pi_3$  totally dominates  $\pi_1$ , which is a contradiction. Hence, there does not exist any policy that weakly totally dominates  $\pi_2$ , meaning that condition b) is certain. We note that  $U_{r_{\pi_1}^*}(\pi_1) = U_{r_{\pi_2}^*}(\pi_1) = U_{r_{\pi_2}^*}(\pi_2)$  and the weak total domination between  $\pi_1, \pi_2$  imply that  $r_{\pi_1}^*, r_{\pi_2}^* \in \arg \max_{r \in R} U_r(\pi_1)$ ,  $r_{\pi_2}^* \in \arg \min_{r \in R} U_r(\pi_2)$ , and thus  $\min_{r \in R} U_r(\pi_2) \leq \max_{\pi \in \Pi_{\neg wtd}} \min_{r \in R} U_r(\pi) = c$ . Again,  $\pi_1 \in \text{MinimaxRegret}(R)$  makes  $\text{Regret}(\pi_1, r_{\pi_1}^*) \leq \text{Regret}(\pi_1, r_{\pi_1}^*) \leq \text{Regret}(\pi, r_{\pi}^*)$  not only hold for  $\pi = \pi_2$  but also for any other policy  $\pi \in \Pi_{\neg wtd}$ , then for any policy  $\pi \in \Pi_{\neg wtd}$ ,  $U_{r_{\pi}^*}(\pi_1) \geq U_{r_{\pi}^*}(\pi) \geq \min_{r \in R} U_r(\pi)$ . Hence,  $U_{r_{\pi}^*}(\pi_1) \geq \max_{\pi \in \Pi_{\neg wtd}} \min_{r \in R} U_r(\pi) = c$ . Since  $U_{r_{\pi}^*}(\pi_1) = \min_{r \in R} U_r(\pi_2)$  as aforementioned,  $\min_{r \in R} U_r(\pi_2) > \max_{\pi \in \Pi_{\neg wtd}} \min_{r \in R} U_r(\pi)$  will cause a contradiction. Hence,  $\min_{r \in R} U_r(\pi_2) = \max_{\pi \in \Pi_{\neg wtd}} \min_{r \in R} U_r(\pi) = c$ . As a result,  $\mathcal{U}_{\mathcal{R}}(\pi) = U_{r_{\pi}^*}(\pi) = \max_{\pi' \in \Pi_{\neg wtd}} \min_{r \in R} U_r(\pi') = c$ , and  $\mathcal{U}_{\mathcal{R}}(\pi) = c - \text{Regret}(\pi, r_{\pi}^*) \geq \max_{\pi' \in \Pi_{\neg wtd}} c - \text{Regret}(\pi', r_{\pi'}^*) = \max_{\pi' \in \Pi_{\neg wtd}} \mathcal{U}_{\mathcal{R}}(\pi')$ . In other words, if a weakly

totally dominated policy  $\pi$  is a solution to  $\text{MinimaxRegret}(R)$ , then its  $\mathcal{U}_{\mathcal{R}}(\pi)$  is no less than that of any non weakly totally dominated policy. This also complete the proof at the end of 1), because if a none weakly totally dominated policy  $\pi_1$  is a solution to  $\arg \max_{\pi \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi)$  but not a solution to  $\text{MinimaxRegret}(R)$ , then  $\text{Regret}(\pi_1, r_{\pi_1}^*) > 0$  and a weakly totally dominated policy  $\pi_2$  must be the solution to  $\text{MinimaxRegret}(R)$ . Then,  $\mathcal{U}_{\mathcal{R}}(\pi_2) = c > c - \text{Regret}(\pi_1, r_{\pi_1}^*) = \mathcal{U}_{\mathcal{R}}(\pi_1)$ , which, however, contradicts  $\pi_1 \in \arg \max_{\pi \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi)$ .

It is obvious that a weakly totally dominated policy  $\pi \in \text{MinimaxRegret}(R)$  has a  $\mathcal{U}_{\mathcal{R}}(\pi)$  no less than any other weakly totally dominated policy. Because for any other weakly totally dominated policy  $\pi_1$ ,  $\mathcal{U}_{\mathcal{R}}(\pi_1) \leq c$  and  $\text{Regret}(\pi_1, r_{\pi_1}^*) \leq \text{Regret}(\pi, r_{\pi}^*)$ , hence  $\mathcal{U}_{\mathcal{R}}(\pi_1) \leq \mathcal{U}_{\mathcal{R}}(\pi)$  according to Eq.13.

So far we have shown that if a weakly totally dominated policy  $\pi$  is a solution to  $\text{MinimaxRegret}(R)$ , it is also a solution to  $\arg \max_{\pi' \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi')$ . Next, we need to show that the reverse is also true, i.e., if a weakly totally dominated policy  $\pi$  is a solution to  $\arg \max_{\pi \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi)$ , it must also be a solution to  $\text{MinimaxRegret}(R)$ . In order to prove its truthfulness, we need to show that if  $\pi \notin \text{MinimaxRegret}(R)$ , whether there exists: a) a none weakly totally dominated policy  $\pi_1$ , or b) another weakly totally dominated policy  $\pi_1$ , such that  $\pi_1 \in \text{MinimaxRegret}(R)$  and  $\mathcal{U}_{\mathcal{R}}(\pi_1) \leq \mathcal{U}_{\mathcal{R}}(\pi)$ . If neither of the two policies exists, we can complete our proof. Since it has been proved in 1) that if a none weakly totally dominated policy achieves  $\text{MinimaxRegret}(R)$ , it also achieves  $\arg \max_{\pi' \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi')$ , the policy described in condition a) does not exist. Hence, it is only necessary to prove that the policy in condition b) also does not exist.

If such weakly totally dominated policy  $\pi_1$  exists,  $\pi \notin \text{MinimaxRegret}(R)$  and  $\pi_1 \in \text{MinimaxRegret}(R)$  indicates  $\text{Regret}(\pi, r_{\pi}^*) > \text{Regret}(\pi_1, r_{\pi_1}^*)$ . Since  $\mathcal{U}_{\mathcal{R}}(\pi_1) \geq \mathcal{U}_{\mathcal{R}}(\pi)$ , according to Eq.13,  $\mathcal{U}_{\mathcal{R}}(\pi_1) = c - \text{Regret}(\pi_1, r_{\pi_1}^*) \leq \mathcal{U}_{\mathcal{R}}(\pi) = \frac{\mathcal{U}_{\mathcal{R}}(\pi) - U_{r_{\pi}^*}(\pi)}{c - U_{r_{\pi}^*}(\pi)} \cdot \mathcal{U}_{\mathcal{R}}(\pi) - \text{Regret}(\pi, r_{\pi}^*)$ . Thus  $\frac{\mathcal{U}_{\mathcal{R}}(\pi) - U_{r_{\pi}^*}(\pi)}{c - U_{r_{\pi}^*}(\pi)}(\pi) \cdot \mathcal{U}_{\mathcal{R}} \geq c + \text{Regret}(\pi, r_{\pi}^*) - \text{Regret}(\pi_1, r_{\pi_1}^*) > c$ , which is impossible due to  $\mathcal{U}_{\mathcal{R}} \leq c$ . Therefore, such  $\pi_1$  also does not exist. In fact, this can be reasoned from another perspective. If there exists a weakly totally dominated policy  $\pi_1$  with  $U_{r_{\pi_1}^*}(\pi_1) = c = U_{r_{\pi}^*}(\pi)$  but  $\pi_1 \notin \text{MinimaxRegret}(R)$ , then  $\text{Regret}(\pi, r_{\pi}^*) > \text{Regret}(\pi_1, r_{\pi_1}^*)$ . It also indicates  $\max_{\pi' \in \Pi} U_{r_{\pi}^*}(\pi') > \max_{\pi' \in \Pi} U_{r_{\pi_1}^*}(\pi')$ . Meanwhile,  $\text{Regret}(\pi_1, r_{\pi}^*) := \max_{\pi' \in \Pi} U_{r_{\pi}^*}(\pi') - U_{r_{\pi}^*}(\pi_1) \leq \text{Regret}(\pi_1, r_{\pi_1}^*) := \max_{\pi' \in \Pi} U_{r_{\pi_1}^*}(\pi') - U_{r_{\pi_1}^*}(\pi_1) := \max_{r \in R} \max_{\pi' \in \Pi} U_r(\pi') - U_r(\pi_1)$  indicates  $\max_{\pi' \in \Pi} U_{r_{\pi}^*}(\pi') - \max_{\pi' \in \Pi} U_{r_{\pi_1}^*}(\pi') \leq U_{r_{\pi}^*}(\pi_1) - U_{r_{\pi_1}^*}(\pi_1)$ . However, we have proved that, for a weakly totally dominated policy,  $\pi_1 \in \text{MinimaxRegret}(R)$  indicates  $U_{r_{\pi_1}^*}(\pi_1) = \max_{r \in R} U_r(\pi_1)$ . Hence,  $\max_{\pi' \in \Pi} U_{r_{\pi}^*}(\pi') - \max_{\pi' \in \Pi} U_{r_{\pi_1}^*}(\pi') \leq U_{r_{\pi}^*}(\pi_1) - U_{r_{\pi_1}^*}(\pi_1) \leq 0$  and it contradicts  $\max_{\pi' \in \Pi} U_{r_{\pi}^*}(\pi') > \max_{\pi' \in \Pi} U_{r_{\pi_1}^*}(\pi')$ . Therefore, such  $\pi_1$  does not exist. In summary, we have exhausted all conditions and can assert that for any policies, being a solution to  $\text{MinimaxRegret}(R)$  is equivalent to a solution to  $\arg \max_{\pi \in \Pi} \mathcal{U}_{\mathcal{R}}(\pi)$ . We complete our proof.  $\square$

### A.3 MEASURING POLICY PERFORMANCE

Recall that the function  $U_r(\pi)$  is used to measure the performance of a policy  $\pi$  under a reward function  $r$ . In Dennis et al. (2020),  $U_r(\pi) = \eta_r(\pi)$ . In this section, we discuss the validity of letting  $U_r(\pi)$  be the loss function of a generic IRL objective, e.g.,  $U_r(\pi) = \eta_r(\pi) - \eta_r(\pi_E)$  where  $\eta_r(\pi_E)$  measures the expected return of the expert policy  $\pi_E$  and can be estimated if an expert demonstration set  $E$  instead of  $\pi_E$  is provided. If further letting  $R_E = \{r | r \in \arg \min_{r' \in R} \max_{\pi \in \Pi} U_{r'}(\pi) - U_{r'}(\pi_E)\}$ ,  $\max_{\pi \in \Pi} U_r(\pi)$  is a constant for any  $r \in R_E$ , notated as  $u := \max_{\pi \in \Pi} U_r(\pi)$ . Because by definition  $R_E = \{r | r \in \arg \min_{r' \in R} \max_{\pi \in \Pi} U_{r'}(\pi)\}$ . If there exists  $r_1, r_2 \in R_E$  such that  $\max_{\pi \in \Pi} U_{r_1}(\pi) < \max_{\pi \in \Pi} U_{r_2}(\pi)$ ,  $r_2$  will not be a member of  $R_E$ . Furthermore,  $\{U_r(\pi) | \pi \in \Pi, r \in R_E\}$  will be upper-bounded by a

constant  $u = \max_{\pi \in \Pi} U_r(\pi)$ . Because if there exists a policy  $\pi \in \Pi$  and a reward function  $r \in R_E$  with  $U_r(\pi) > u$ , it contradicts the fact that  $u = \max_{\pi' \in \Pi} U_r(\pi')$ . In this case,  $\text{MinimaxRegret}(R_E) = \arg \min_{\pi \in \Pi} \max_{r \in R_E} \text{Regret}(\pi, r) = \arg \min_{\pi \in \Pi} \max_{r \in R_E} u - U_r(\pi) = \arg \max_{\pi \in \Pi} \min_{r \in R_E} U_r(\pi)$ . Note that before making any other assumption on  $R_E$ ,  $\Pi$  and  $U_r(\cdot)$ ,  $\max_{\pi \in \Pi} \min_{r \in R_E} U_r(\pi)$  cannot be regarded as the same as IRL itself  $\min_{r \in R} \max_{\pi \in \Pi} U_r(\pi)$ . The solution to  $\arg \max_{\pi \in \Pi} \min_{r \in R_E} U_r(\pi)$  is the policy with the highest worst case  $U_r(\pi)$  for  $r \in R_E$ . The IRL problem however may induce a policy that maximizes  $U_r(\cdot)$  for some  $r \in R_E$  while minimizing  $U_{r'}(\cdot)$  for some other  $r' \in R_E$ . While  $\min_{r \in R} \max_{\pi \in \Pi} U_r(\pi) = u$ , it is possible that  $\max_{\pi \in \Pi} \min_{r \in R_E} U_r(\pi) < u$ . In fact, it is easily observable that the solutions to  $\text{MinimaxRegret}$  with some  $U_r(\pi)$  will be the same as that of letting  $U_r(\pi) := U_r(\pi) - U_r(\pi_E)$ . Hence, in this paper we simply use  $\eta_r(\pi)$  as  $U_r(\pi)(R)$ .

**Lemma 5.** *If a policy  $\pi \in \text{MinimaxRegret}(R_E)$  when the policy performance is measured with some  $U_r$ , then  $\pi \in \text{MinimaxRegret}(R_E)$  when letting  $U_r(\pi) := U_r(\pi) - U_r(\pi_E)$ .*

*Proof.* When using  $U_r(\pi) := U_r(\pi) - U_r(\pi_E)$  to measure the policy performance, solving  $\text{MinimaxRegret}(R)$  is to solve Eq. [14](#), which is the same as Eq. [9](#).

$$\begin{aligned} \text{MinimaxRegret}(R_E) &= \arg \max_{\pi \in \Pi} \min_{r \in R_E} \text{Regret}(\pi, r) \\ &= \arg \max_{\pi \in \Pi} \min_{r \in R_E} \max_{\pi' \in \Pi} \{U_r(\pi') - U_r(\pi_E)\} - (U_r(\pi) - U_r(\pi_E)) \\ &= \arg \max_{\pi \in \Pi} \min_{r \in R_E} \max_{\pi' \in \Pi} U_r(\pi') - U_r(\pi) \end{aligned} \quad (14)$$

□

#### A.4 CRITERION FOR SUCCESSFUL POLICY LEARNING

**Theorem [1](#)** (Task-Failure Avoidance) If the following conditions (1) (2) hold for  $R$ , then the optimal protagonist policy  $\pi_P := \text{MinimaxRegret}(R)$  satisfies that  $\forall r_{al} \in R, U_{r_{al}}(\pi_P) \notin F_{r_{al}}$ .

- (1) There exists  $r_{al} \in R$ , and  $\max_{r_{al} \in R} \{\sup F_{r_{al}} - \inf F_{r_{al}}\} < \min_{r_{al} \in R} \{\inf S_{r_{al}} - \sup F_{r_{al}}\} \wedge \max_{r_{al} \in R} \{\sup S_{r_{al}} - \inf S_{r_{al}}\} < \min_{r_{al} \in R} \{\inf S_{r_{al}} - \sup F_{r_{al}}\}$ ;
- (2) There exists a policy  $\pi^*$  such that  $\forall r_{al} \in R, U_{r_{al}}(\pi^*) \in S_{r_{al}}$ , and  $\forall r_{mis} \in R, \max_{\pi \in \Pi} U_{r_{mis}}(\pi) - U_{r_{mis}}(\pi^*) < \min_{r_{al} \in R} \{\inf S_{r_{al}} - \sup F_{r_{al}}\}$ .

*Proof.* Suppose the conditions are met, and a policy  $\pi_1$  satisfies the property described in conditions 2). Then for any policy  $\pi_2 \in \text{MinimaxRegret}(R)$ , if  $\pi_2$  does not satisfy the mentioned property, there exists a task-aligned reward function  $r_{al} \in R_E$  such that  $U_{r_{al}}(\pi_2) \in F_{r_{al}}$ . In this case  $\text{Regret}(\pi_2, r_{al}) = \max_{\pi \in \Pi} U_{r_{al}}(\pi) - U_{r_{al}}(\pi_2) \geq \inf S_{r_{al}} - \sup F_{r_{al}} \geq \min_{r'_{al} \in R_E} \{\inf S_{r'_{al}} - \sup F_{r'_{al}}\}$ . However, for  $\pi_1$ , it holds for any task-aligned reward function  $\hat{r}_{al} \in R_E$  that  $\text{Regret}(\pi_2, \hat{r}_{al}) \leq \sup S_{\hat{r}_{al}} - \inf S_{\hat{r}_{al}} < \min_{r'_{al} \in R_E} \{\inf S_{r'_{al}} - \sup F_{r'_{al}}\}$ , and it also holds for any misaligned reward function  $r_{mis} \in R_E$  that  $\text{Regret}(\pi_2, r_{mis}) = \max_{\pi \in \Pi} U_{r_{mis}}(\pi) - U_{r_{mis}}(\pi_2) < \min_{r'_{al} \in R_E} \{\inf S_{r_{al}} - \sup F_{r_{al}}\}$ . Hence,  $\text{Regret}(\pi_2, r_{al}) < \text{Regret}(\pi_1, r_{al})$ , contradicting  $\pi_1 \in \text{MiniRegret}$ . We complete the proof. □

**Corollary [1](#)** (Task-Success Guarantee) Assume that Condition (1) in Theorem [1](#) is satisfied. If there exists a policy  $\pi^*$  such that  $\forall r \in R, \max_{\pi} U_r(\pi) - U_r(\pi^*) < \min_{r_{al} \in R} \{\sup S_{r_{al}} - \inf S_{r_{al}}\}$ , then the optimal protagonist policy  $\pi_P := \text{MinimaxRegret}(R)$  satisfies that  $\forall r_{al} \in R, U_{r_{al}}(\pi) \in S_{r_{al}}$ .



*Proof.* Since  $\max_{r \in R} \max_{\pi} U_r(\pi) - U_r(\pi_P) \leq \max_{r \in R} \max_{\pi} U_r(\pi) - U_r(\pi^*) < \min_{r_{al} \in R} \{\sup S_{r_{al}} - \inf S_{r_{al}}\}$ , we can conclude that for any  $r_{al} \in R$ ,  $\max_{\pi} U_{r_{al}}(\pi) - U_{r_{al}}(\pi_P) \leq \{\sup S_{r_{al}} - \inf S_{r_{al}}\}$ , in other words,  $U_{r_{al}}(\pi_P) \in S_{r_{al}}$ . The proof is complete.  $\square$

**Theorem 2**(Task-Failure Avoidance) If the following conditions (1) (2) hold for  $R_{E,\delta}$ , then the optimal protagonist policy  $\pi_P := \text{MinimaxRegret}(R_{E,\delta})$  satisfies that  $\forall r_{al} \in R_{E,\delta}$ ,  $U_{r_{al}}(\pi) \notin F_{r_{al}}$ .

- (1) The condition (1) in Theorem 1 holds
- (2)  $\forall r_{al} \in R_{E,\delta}$ ,  $L_{r_{al}} \cdot W_E - \delta \leq \sup S_{r_{al}} - \inf S_{r_{al}}$  and  $\forall r_{mis} \in R_{E,\delta}$ ,  $L_{r_{mis}} \cdot W_E - \delta < \min_{r_{al} \in R_{E,\delta}} \{\inf S_{r_{al}} - \sup F_{r_{al}}\}$ .

*Proof.* Without loss of generality, we consider  $U_r(\pi) = \mathbb{E}_{\tau \sim \pi}[r(\tau)]$ . Since  $W_E \triangleq \min_{\pi \in \Pi} W_1(\pi, E) = \frac{1}{K} \sup_{\|r\|_L \leq K} U_r(E) - U_r(\pi)$  for any  $K > 0$ , let  $\pi^*$  be the policy that achieves the minimality in  $W_E$ . Then for any  $r_{al} \in R$ , the term  $L_{r_{al}} \cdot W_E - \delta \geq L_{r_{al}} \cdot \frac{1}{L_{r_{al}}} \sup_{\|r\|_L \leq L_{r_{al}}} U_r(E) - U_r(\pi) \geq U_{r_{al}}(E) - U_{r_{al}}(\pi)$ . Hence, for all  $r_{al} \in R$ ,  $U_{r_{al}}(E) - U_{r_{al}}(\pi) < \sup S_{r_{al}} - \inf S_{r_{al}}$ , i.e.,  $U_{r_{al}}(\pi^*) \in S_{r_{al}}$ . Likewise,  $L_{r_{mis}} \cdot W_E - \delta < \min_{r_{al} \in R_{E,\delta}} \{\inf S_{r_{al}} - \sup F_{r_{al}}\}$  indicates that for all  $r_{mis} \in R$ ,  $U_{r_{mis}}(E) - U_{r_{mis}}(\pi) < \min_{r_{al} \in R_{E,\delta}} \{\inf S_{r_{al}} - \sup F_{r_{al}}\}$ . Then, we have recovered the condition (2) in Theorem 1. As a result, we deliver the same guarantees in Theorem 1.  $\square$

**Corollary 2**(Task-Success Guarantee) Assume that the condition (1) in Theorem 1 holds for  $R_{E,\delta}$ . If for any  $r \in R_{E,\delta}$ ,  $L_r \cdot W_E - \delta \leq \min_{r_{al} \in R} \{\sup S_{r_{al}} - \inf S_{r_{al}}\}$ , then the optimal protagonist policy  $\pi_P = \text{MinimaxRegret}(R_{E,\delta})$  satisfies  $\forall r_{al} \in R_{E,\delta}$ ,  $U_{r_{al}}(\pi) \in S_{r_{al}}$ .

*Proof.* Again, we let  $\pi^*$  be the policy that achieves the minimality in  $W_E$ . Then, we have  $L_r \cdot W_E - \delta \geq L_r \cdot \frac{1}{L_r} \sup_{\|r\|_L \leq L_r} U_r(E) - U_r(\pi^*) \geq U_r(E) - U_r(\pi^*)$  for any  $r \in R$ . We have recovered the condition in Corollary 1. The proof is complete.  $\square$

**Theorem 3** Let the IRL loss be in the form of  $J_{IRL}(\pi, r) := U_r(\pi) - U_r(\pi_E)$  for some  $U_r(\pi)$ . If  $\arg \min_{r \in R} \max_{\pi \in \Pi} J_{IRL}(r_E, \pi)$  can reach Nash Equilibrium with a reward function set  $R_E$  and a policy set  $\Pi_E$ , then  $\Pi_E$  equals the set of solutions to *MinimaxRegret*.

*Proof.* The reward function set  $R_E$  and the policy set  $\Pi_E$  achieving Nash Equilibrium for  $\arg \min_{r \in R} \max_{\pi \in \Pi} J_{IRL}(r_E, \pi)$  indicates that for any  $r \in R_E$ ,  $\pi \in \Pi_E$ ,  $\pi \in \arg \max_{\pi \in \Pi} U_r(\pi) - U_r(\pi_E)$ .

Then  $\Pi_E$  will be the solution to  $\arg \max_{\pi_P \in \Pi} \min_{r \in R_E} \left\{ \max_{\pi_A \in \Pi} U_r(\pi_A) - U_r(\pi_E) \right\} - (U_r(\pi_P) - U_r(\pi_E))$  because the policies in  $\Pi_E$  achieve zero regret. Then Lemma 5 states that  $\Pi_E$  will also be the solution to  $\arg \max_{\pi_P \in \Pi} \min_{r \in R_E} \left\{ \max_{\pi_A \in \Pi} U_r(\pi_A) \right\} - U_r(\pi_P)$ . We finish the proof.  $\square$

## A.5 RUNNING EXAMPLE

**Example 1** In a two-state transition system where  $s_0$  is the initial state, and  $s_1$  is an absorbing state, an agent can choose action  $a_0$  at state  $s_0$  to stay at  $s_0$  or choose  $a_1$  to reach  $s_1$ . Any agent can start from  $s_0$  and choose actions for 5 timesteps. The task is to learn a stochastic policy to reach  $s_1$ . Expert only demonstrates one trajectory  $E = \{\tau_E = (s_0, a_1, s_1, s_1, s_1)\}$ , i.e., choose  $a_1$  at  $s_0$  and then stay in  $s_1$  for the rest 4 steps. The convex combinations of two basis functions  $r_1$  and  $r_2$  constitute the hypothesis set of reward functions  $R = \{r | r(s, a) = \alpha \cdot r_1(s, a) + (1 - \alpha) \cdot r_2(s, a), \alpha \in [0, 1]\}$  where  $r_1(s, a) \equiv 4$  constantly, and  $r_2(s_0, a_0) = r_2(s_0, a_1) = 0, r_2(s_1, \cdot) = 5$ . We now prove that when applying IRL to learn a reward function from  $E$ , any convex combination of  $r_1$  and  $r_2$  is an

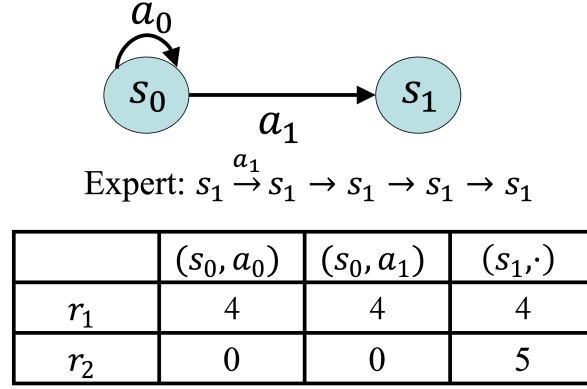


Figure 5: ss

optimal solution, i.e.,  $R_E = R$ . Furthermore, we prove that  $\text{MinimaxRegret}(R)$  will produce  $\pi_P(a_1|s_0) = 1$  which is the desired solution.

*Proof.* For any policy  $\pi$  and reward function  $r \in R$ , the policy performance  $U_r(\pi)$  can be represented as follows.

$$\begin{aligned}
U_r(\pi) &= 4\alpha\pi(a_0|s_0)^t + \sum_{t=0}^t (4\alpha t + 4\alpha + (4\alpha + 5(1-\alpha))(4-t)) \cdot \pi(a_0|s_0)^t (1 - \pi(a_0|s_0)) \\
&= 20 - 5(1-\alpha) \sum_{t=1}^4 \pi(a_0|s_0)^4
\end{aligned} \tag{15}$$

Then for any  $\alpha \in [0, 1]$ ,  $\max_{\pi} U_r(\pi) = 20$ , and the optimal policy is  $\pi^*(a_0|s_0) = 0$ ,  $\pi^*(a_1|s_0) = 1$ . For  $\alpha = 1$ ,  $U_r(\pi) = 20$  constantly for any  $\pi$ . For expert demonstration  $E$ ,  $U_r(E) = 4\alpha \cdot 5 + 5(1-\alpha) \cdot 4 \equiv 20$ . Therefore, any  $r = \alpha r_1 + (1-\alpha)r_2$  with  $\alpha \in [0, 1]$  can be the optimal solution to IRL since they all induce the same performance gap  $U_r(E) - \max_{\pi} U_r(\pi) = 0$ . This ambiguity in choosing reward functions can cause reward misalignment because the reward function  $r = r_1$  induced by  $\alpha = 1$  violates the criterion of task-reward alignment in Definition 1 achieving optimality under  $r_1$  does not guarantee task success since all policies are optimal under  $r_1$ . On the other hand, it is obvious that  $\text{MinimaxRegret}(R)$  induces a policy  $\pi_P(a_1|s_0) = 1$ . Because this policy has the lowest worst-case regret  $\text{Regret}(\pi_P, r) = \max_{r \in R} \max_{\pi_A} U_r(\pi_A) - U_r(\pi_P) = 0$ .  $\square$

## B APPROACH TO SOLVING MINIMAXREGRET

In this section, we develop a series of theories that lead to two bounds of the Protagonist Antagonist Induced Regret. By using those bounds, we formulate objective functions for solving Imitation Learning problems with PAGAR.

### B.1 PROTAGONIST ANTAGONIST INDUCED REGRET BOUNDS

Our theories are inspired by the on-policy policy improvement methods in Schulman et al. (2015). The theories in Schulman et al. (2015) are under the setting where entropy regularizer is not considered. In our implementation, we always consider entropy regularized RL of which the objective is to learn a policy that maximizes  $J_{RL}(\pi; r) = \eta_r(\pi) + \mathcal{H}(\pi)$ . Also, since we use GAN-based IRL algorithms, the learned reward function  $r$  as proved by Fu et al. (2018) is a distribution. Moreover, it is also proved in Fu et al. (2018) that a policy  $\pi$  being optimal under  $r$  indicates that  $\log \pi \equiv r \equiv \mathcal{A}_{\pi}$ . We omit the proof and let the reader refer to Fu et al. (2018) for details. Although all our theories are about the relationship between the Protagonist Antagonist Induced Regret and the soft advantage function  $\mathcal{A}_{\pi}$ , the equivalence between  $\mathcal{A}_{\pi}$  and  $r$  allows us to use the theories to formulate our



reward optimization objective functions. To start off, we denote the reward function to be optimized as  $r$ . Given the intermediate learned reward function  $r$ , we study the Protagonist Antagonist Induced Regret between two policies  $\pi_1$  and  $\pi_2$ .

**Lemma 6.** *Given a reward function  $r$  and a pair of policies  $\pi_1$  and  $\pi_2$ ,*

$$\eta_r(\pi_1) - \eta_r(\pi_2) = \mathbb{E}_{\tau \sim \pi_1} \left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)}) \right] + \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{H}(\pi_2(\cdot | s^{(t)})) \right] \quad (16)$$

*Proof.* This proof follows the proof of Lemma 1 in [Schulman et al. \(2015\)](#) where RL is not entropy-regularized. For entropy-regularized RL, since  $\mathcal{A}_{\pi}(s, a^{(t)}) = \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a^{(t)})} [r(s, a^{(t)}) + \gamma \mathcal{V}_{\pi}(s') - \mathcal{V}_{\pi}(s)]$ ,

$$\begin{aligned} & \mathbb{E}_{\tau \sim \pi_1} \left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)}) \right] \\ &= \mathbb{E}_{\tau \sim \pi_1} \left[ \sum_{t=0}^{\infty} \gamma^t \left( r(s^{(t+1)}, a^{(t+1)}) + \gamma \mathcal{V}_{\pi_2}(s^{(t+1)}) - \mathcal{V}_{\pi_2}(s^{(t)}) \right) \right] \\ &= \mathbb{E}_{\tau \sim \pi_1} \left[ \sum_{t=0}^{\infty} \gamma^t r(s^{(t)}, a^{(t)}) - \mathcal{V}_{\pi_2}(s^{(0)}) \right] \\ &= \mathbb{E}_{\tau \sim \pi_1} \left[ \sum_{t=0}^{\infty} \gamma^t r(s^{(t)}, a^{(t)}) \right] - \mathbb{E}_{s^{(0)} \sim d_0} [\mathcal{V}_{\pi_2}(s^{(0)})] \\ &= \mathbb{E}_{\tau \sim \pi_1} \left[ \sum_{t=0}^{\infty} \gamma^t r(s^{(t)}, a^{(t)}) \right] - \mathbb{E}_{\tau \sim \pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t r(s^{(t)}, a^{(t)}) + \mathcal{H}(\pi_2(\cdot | s^{(t)})) \right] \\ &= \eta_r(\pi_1) - \eta_r(\pi_2) - \mathbb{E}_{\tau \sim \pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{H}(\pi_2(\cdot | s^{(t)})) \right] \\ &= \eta_r(\pi_1) - \eta_r(\pi_2) - \mathcal{H}(\pi_2) \end{aligned}$$

□

**Remark 1.** Lemma 6 confirms that  $\mathbb{E}_{\tau \sim \pi} [\sum_{t=0}^{\infty} \gamma^t \mathcal{A}_{\pi}(s^{(t)}, a^{(t)})] = \eta_r(\pi) - \eta_r(\pi) + \mathcal{H}(\pi) = \mathcal{H}(\pi)$ .

We follow [Schulman et al. \(2015\)](#) and denote  $\Delta \mathcal{A}(s) = \mathbb{E}_{a \sim \pi_1(\cdot | s)} [\mathcal{A}_{\pi_2}(s, a)] - \mathbb{E}_{a \sim \pi_2(\cdot | s)} [\mathcal{A}_{\pi_2}(s, a)]$  as the difference between the expected advantages of following  $\pi_2$  after choosing an action respectively by following policy  $\pi_1$  and  $\pi_2$  at any state  $s$ . Although the setting of [Schulman et al. \(2015\)](#) differs from ours by having the expected advantage  $\mathbb{E}_{a \sim \pi_2(\cdot | s)} [\mathcal{A}_{\pi_2}(s, a)]$  equal to 0 due to the absence of entropy regularization, the following definition and lemmas from [Schulman et al. \(2015\)](#) remain valid in our setting.

**Definition 7.** [Schulman et al. \(2015\)](#), the protagonist policy  $\pi_1$  and the antagonist policy  $\pi_2$  are  $\alpha$ -coupled if they defines a joint distribution over  $(a, \tilde{a}) \in \mathbb{A} \times \mathbb{A}$ , such that  $\text{Prob}(a \neq \tilde{a} | s) \leq \alpha$  for all  $s$ .

**Lemma 7.** [Schulman et al. \(2015\)](#) Given that the protagonist policy  $\pi_1$  and the antagonist policy  $\pi_2$  are  $\alpha$ -coupled, then for all state  $s$ ,

$$|\Delta \mathcal{A}(s)| \leq 2\alpha \max_a |\mathcal{A}_{\pi_2}(s, a)| \quad (17)$$

**Lemma 8.** [Schulman et al. \(2015\)](#) Given that the protagonist policy  $\pi_1$  and the antagonist policy  $\pi_2$  are  $\alpha$ -coupled, then

$$\left| \mathbb{E}_{s^{(t)} \sim \pi_1} [\Delta \mathcal{A}(s^{(t)})] - \mathbb{E}_{s^{(t)} \sim \pi_2} [\Delta \mathcal{A}(s^{(t)})] \right| \leq 4\alpha(1 - (1 - \alpha)^t) \max_{s, a} |\mathcal{A}_{\pi_2}(s, a)| \quad (18)$$

**Lemma 9.** *Given that the protagonist policy  $\pi_1$  and the antagonist policy  $\pi_2$  are  $\alpha$ -coupled, then*

$$\mathbb{E}_{\substack{s^{(t)} \sim \pi_1 \\ a^{(t)} \sim \pi_2}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] - \mathbb{E}_{\substack{s^{(t)} \sim \pi_2 \\ a^{(t)} \sim \pi_2}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \leq 2(1 - (1 - \alpha)^t) \max_{(s,a)} |\mathcal{A}_{\pi_2}(s, a)| \quad (19)$$

*Proof.* The proof is similar to that of Lemma 8 in [Schulman et al. \(2015\)](#). Let  $n_t$  be the number of times that  $a^{(t')} \sim \pi_1$  does not equal  $a^{(t')} \sim \pi_2$  for  $t' < t$ , i.e., the number of times that  $\pi_1$  and  $\pi_2$  disagree before timestep  $t$ . Then for  $s^{(t)} \sim \pi_1$ , we have the following.

$$\begin{aligned} & \mathbb{E}_{\substack{s^{(t)} \sim \pi_1 \\ a^{(t)} \sim \pi_2}} [\mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})]] \\ = & P(n_t = 0) \mathbb{E}_{\substack{s^{(t)} \sim \pi_1 \\ n_t = 0}} [\mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})]] + P(n_t > 0) \mathbb{E}_{\substack{s^{(t)} \sim \pi_1 \\ n_t > 0}} [\mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})]] \end{aligned}$$

The expectation decomposes similarly for  $s^{(t)} \sim \pi_2$ .

$$\begin{aligned} & \mathbb{E}_{\substack{s^{(t)} \sim \pi_2 \\ a^{(t)} \sim \pi_2}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \\ = & P(n_t = 0) \mathbb{E}_{\substack{s^{(t)} \sim \pi_2 \\ a^{(t)} \sim \pi_2 \\ n_t = 0}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] + P(n_t > 0) \mathbb{E}_{\substack{s^{(t)} \sim \pi_2 \\ a^{(t)} \sim \pi_2 \\ n_t > 0}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \end{aligned}$$

When computing  $\mathbb{E}_{\substack{s^{(t)} \sim \pi_1 \\ a^{(t)} \sim \pi_2}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] - \mathbb{E}_{\substack{s^{(t)} \sim \pi_2 \\ a^{(t)} \sim \pi_2}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})]$ , the terms with  $n_t =$

0 cancel each other because  $n_t = 0$  indicates that  $\pi_1$  and  $\pi_2$  agreed on all timesteps less than  $t$ . That leads to the following.

$$\begin{aligned} & \mathbb{E}_{\substack{s^{(t)} \sim \pi_1 \\ a^{(t)} \sim \pi_2}} [\mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})]] - \mathbb{E}_{\substack{s^{(t)} \sim \pi_2 \\ a^{(t)} \sim \pi_2}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \\ = & P(n_t > 0) \mathbb{E}_{\substack{s^{(t)} \sim \pi_1 \\ n_t > 0}} [\mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})]] - P(n_t > 0) \mathbb{E}_{\substack{s^{(t)} \sim \pi_2 \\ a^{(t)} \sim \pi_2 \\ n_t > 0}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \end{aligned}$$

By definition of  $\alpha$ , the probability of  $\pi_1$  and  $\pi_2$  agreeing at timestep  $t'$  is no less than  $1 - \alpha$ . Hence,  $P(n_t > 0) \leq 1 - (1 - \alpha)^t$ . Hence, we have the following bound.

$$\begin{aligned} & \left| \mathbb{E}_{\substack{s^{(t)} \sim \pi_1 \\ a^{(t)} \sim \pi_2}} [\mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})]] - \mathbb{E}_{\substack{s^{(t)} \sim \pi_2 \\ a^{(t)} \sim \pi_2}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right| \\ = & \left| P(n_t > 0) \mathbb{E}_{\substack{s^{(t)} \sim \pi_1 \\ n_t > 0}} [\mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})]] - P(n_t > 0) \mathbb{E}_{\substack{s^{(t)} \sim \pi_2 \\ a^{(t)} \sim \pi_2 \\ n_t > 0}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right| \\ \leq & P(n_t > 0) \left( \left| \mathbb{E}_{\substack{s^{(t)} \sim \pi_1 \\ a^{(t)} \sim \pi_2 \\ n_t \geq 0}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right| + \left| \mathbb{E}_{\substack{s^{(t)} \sim \pi_2 \\ a^{(t)} \sim \pi_2 \\ n_t > 0}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right| \right) \\ \leq & 2(1 - (1 - \alpha)^t) \max_{(s,a)} |\mathcal{A}_{\pi_2}(s, a)| \quad (20) \end{aligned}$$

□

The preceding lemmas lead to the proof for Theorem 4 in the main text.

**Theorem 4.** Suppose that  $\pi_2$  is the optimal policy in terms of entropy regularized RL under  $r$ . Let  $\alpha = \max_s D_{TV}(\pi_1(\cdot|s), \pi_2(\cdot|s))$ ,  $\epsilon = \max_{s,a} |\mathcal{A}_{\pi_2}(s, a^{(t)})|$ , and  $\Delta\mathcal{A}(s) = \mathbb{E}_{a \sim \pi_1} [\mathcal{A}_{\pi_2}(s, a)] - \mathbb{E}_{a \sim \pi_2} [\mathcal{A}_{\pi_2}(s, a)]$ . For any policy  $\pi_1$ , the following bounds hold.

$$\left| \eta_r(\pi_1) - \eta_r(\pi_2) - \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_1} [\Delta\mathcal{A}(s^{(t)})] \right| \leq \frac{2\alpha\gamma\epsilon}{(1-\gamma)^2} \quad (21)$$

$$\left| \eta_r(\pi_1) - \eta_r(\pi_2) - \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_2} [\Delta\mathcal{A}(s^{(t)})] \right| \leq \frac{2\alpha\gamma(2\alpha+1)\epsilon}{(1-\gamma)^2} \quad (22)$$

*Proof.* We first leverage Lemma 6 to derive Eq 23. Note that since  $\pi_2$  is optimal under  $r$ , Remark 1 confirmed that  $\mathcal{H}(\pi_2) = -\sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_2} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right]$ .

$$\begin{aligned} & \eta_r(\pi_1) - \eta_r(\pi_2) \\ &= (\eta_r(\pi_1) - \eta_r(\pi_2) - \mathcal{H}(\pi_2)) + \mathcal{H}(\pi_2) \\ &= \mathbb{E}_{\tau \sim \pi_1} \left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)}) \right] + \mathcal{H}(\pi_2) \\ &= \mathbb{E}_{\tau \sim \pi_1} \left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)}) \right] - \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_2} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] \\ &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_1} \left[ \mathbb{E}_{a^{(t)} \sim \pi_1} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] - \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] + \\ & \quad \sum_{t=0}^{\infty} \gamma^t \left( \mathbb{E}_{s^{(t)} \sim \pi_1} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] - \mathbb{E}_{s^{(t)} \sim \pi_2} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] \right) \\ &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_1} [\Delta\mathcal{A}(s^{(t)})] + \\ & \quad \sum_{t=0}^{\infty} \gamma^t \left( \mathbb{E}_{s^{(t)} \sim \pi_1} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] - \mathbb{E}_{s^{(t)} \sim \pi_2} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] \right) \quad (23) \end{aligned}$$

We switch terms between Eq 23 and  $\eta_r(\pi_1) - \eta_r(\pi_2)$ , then use Lemma 9 to derive Eq 24.

$$\begin{aligned} & \left| \eta_r(\pi_1) - \eta_r(\pi_2) - \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_1} [\Delta\mathcal{A}(s^{(t)})] \right| \\ &= \left| \sum_{t=0}^{\infty} \gamma^t \left( \mathbb{E}_{s^{(t)} \sim \pi_1} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] - \mathbb{E}_{s^{(t)} \sim \pi_2} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] \right) \right| \\ &\leq \sum_{t=0}^{\infty} \gamma^t \cdot 2 \max_{(s,a)} |\mathcal{A}_{\pi_2}(s, a)| \cdot (1 - (1 - \alpha)^t) \leq \frac{2\alpha\gamma \max_{(s,a)} |\mathcal{A}_{\pi_2}(s, a)|}{(1-\gamma)^2} \quad (24) \end{aligned}$$

Alternatively, we can expand  $\eta_r(\pi_2) - \eta_r(\pi_1)$  into Eq.25. During the process,  $\mathcal{H}(\pi_2)$  is converted into  $-\sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_2} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right]$ .

$$\begin{aligned}
& \eta_r(\pi_1) - \eta_r(\pi_2) \\
&= (\eta_r(\pi_1) - \eta_r(\pi_2) - \mathcal{H}(\pi_2)) + \mathcal{H}(\pi_2) \\
&= \mathbb{E}_{\tau \sim \pi_1} \left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)}) \right] + \mathcal{H}(\pi_2) \\
&= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_1} \left[ \mathbb{E}_{a^{(t)} \sim \pi_1} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] + \mathcal{H}(\pi_2) \\
&= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_1} \left[ \Delta \mathcal{A}(s^{(t)}) + \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] + \mathcal{H}(\pi_2) \\
&= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_2} \left[ \mathbb{E}_{a^{(t)} \sim \pi_1} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] - \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] - \Delta \mathcal{A}(s^{(t)}) \right] + \\
&\quad \mathbb{E}_{s^{(t)} \sim \pi_1} \left[ \Delta \mathcal{A}(s^{(t)}) + \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] - \mathbb{E}_{\substack{s^{(t)} \sim \pi_2 \\ a^{(t)} \sim \pi_2}} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \\
&= \sum_{t=0}^{\infty} \gamma^t \left( \mathbb{E}_{s^{(t)} \sim \pi_1} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] - 2 \mathbb{E}_{s^{(t)} \sim \pi_2} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] \right) + \\
&\quad \sum_{t=0}^{\infty} \gamma^t \left( \mathbb{E}_{s^{(t)} \sim \pi_2} \left[ \mathbb{E}_{a^{(t)} \sim \pi_1} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] - \left( \mathbb{E}_{s^{(t)} \sim \pi_2} [\Delta \mathcal{A}(s^{(t)})] - \mathbb{E}_{s^{(t)} \sim \pi_1} [\Delta \mathcal{A}(s^{(t)})] \right) \right)
\end{aligned} \tag{25}$$

We switch terms between Eq.25 and  $\eta_r(\pi_1) - \eta_r(\pi_2)$ , then base on Lemma 8 and 9 to derive the inequality in Eq.26.

$$\begin{aligned}
& \left| \eta_r(\pi_1) - \eta_r(\pi_2) - \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_2} [\Delta \mathcal{A}(s^{(t)})] \right| \\
&= \left| \eta_r(\pi_1) - \eta_r(\pi_2) - \right. \\
&\quad \left. \sum_{t=0}^{\infty} \gamma^t \left( \mathbb{E}_{s^{(t)} \sim \pi_2} \left[ \mathbb{E}_{a^{(t)} \sim \pi_1} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] - \mathbb{E}_{s^{(t)} \sim \pi_2} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] \right) \right| \\
&= \left| \sum_{t=0}^{\infty} \gamma^t \left( \mathbb{E}_{s^{(t)} \sim \pi_2} [\Delta \mathcal{A}(s^{(t)})] - \mathbb{E}_{s^{(t)} \sim \pi_1} [\Delta \mathcal{A}(s^{(t)})] \right) - \right. \\
&\quad \left. \sum_{t=0}^{\infty} \gamma^t \left( \mathbb{E}_{s^{(t)} \sim \pi_1} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] - \mathbb{E}_{s^{(t)} \sim \pi_2} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] \right) \right| \\
&\leq \left| \sum_{t=0}^{\infty} \gamma^t \left( \mathbb{E}_{s^{(t)} \sim \pi_2} [\Delta \mathcal{A}(s^{(t)})] - \mathbb{E}_{s^{(t)} \sim \pi_1} [\Delta \mathcal{A}(s^{(t)})] \right) \right| + \\
&\quad \left| \sum_{t=0}^{\infty} \gamma^t \left( \mathbb{E}_{s^{(t)} \sim \pi_1} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] - \mathbb{E}_{s^{(t)} \sim \pi_2} \left[ \mathbb{E}_{a^{(t)} \sim \pi_2} [\mathcal{A}_{\pi_2}(s^{(t)}, a^{(t)})] \right] \right) \right| \\
&\leq \sum_{t=0}^{\infty} \gamma^t \left( (1 - (1 - \alpha)^t) (4\alpha \max_{s,a} |\mathcal{A}_{\pi_2}(s, a)| + 2 \max_{(s,a)} |\mathcal{A}_{\pi_2}(s, a)|) \right) \\
&\quad \frac{2\alpha\gamma(2\alpha + 1) \max_{s,a} |\mathcal{A}_{\pi_2}(s, a)|}{(1 - \gamma)^2}
\end{aligned} \tag{26}$$

It is stated in [Schulman et al. \(2015\)](#) that  $\max D_{TV}(\pi_2(\cdot|s), \pi_1(\cdot|s)) \leq \alpha$ . Hence, by letting  $\alpha := \max_s D_{TV}(\pi_2(\cdot|s), \pi_1(\cdot|s))$ , Eq. 23 and 26 still hold. Then, we have proved Theorem 4.  $\square$

## B.2 OBJECTIVE FUNCTIONS OF REWARD OPTIMIZATION

To derive  $J_{R,1}$  and  $J_{R,2}$ , we let  $\pi_1 = \pi_P$  and  $\pi_2 = \pi_A$ . Then based on Eq. 21 and 22 we derive the following upper-bounds of  $\eta_r(\pi_P) - \eta_r(\pi_A)$ .

$$\eta_r(\pi_P) - \eta_r(\pi_A) \leq \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_P} [\Delta \mathcal{A}(s^{(t)})] + \frac{2\alpha\gamma(2\alpha+1)\epsilon}{(1-\gamma)^2} \quad (27)$$

$$\eta_r(\pi_P) - \eta_r(\pi_A) \geq \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_A} [\Delta \mathcal{A}(s^{(t)})] - \frac{2\alpha\gamma\epsilon}{(1-\gamma)^2} \quad (28)$$

By our assumption that  $\pi_A$  is optimal under  $r$ , we have  $\mathcal{A}_{\pi_A} \equiv r$  [Fu et al. \(2018\)](#). This equivalence enables us to replace  $\mathcal{A}_{\pi_A}$ 's in  $\Delta \mathcal{A}$  with  $r$ . As for the  $\frac{2\alpha\gamma(2\alpha+1)\epsilon}{(1-\gamma)^2}$  and  $\frac{2\alpha\gamma\epsilon}{(1-\gamma)^2}$  terms, since the objective is to maximize  $\eta_r(\pi_A) - \eta_r(\pi_B)$ , we heuristically estimate the  $\epsilon$  in Eq. 27 by using the samples from  $\pi_P$  and the  $\epsilon$  in Eq. 28 by using the samples from  $\pi_A$ . As a result we have the objective functions defined as Eq. 29 and 30 where  $\delta_1(s, a) = \frac{\pi_P(a^{(t)}|s^{(t)})}{\pi_A(a^{(t)}|s^{(t)})}$  and  $\delta_2 = \frac{\pi_A(a^{(t)}|s^{(t)})}{\pi_P(a^{(t)}|s^{(t)})}$  are the importance sampling probability ratio derived from the definition of  $\Delta \mathcal{A}$ ;  $C_1 \propto -\frac{\gamma\hat{\alpha}}{(1-\gamma)}$  and  $C_2 \propto \frac{\gamma\hat{\alpha}}{(1-\gamma)}$  where  $\hat{\alpha}$  is either an estimated maximal KL-divergence between  $\pi_A$  and  $\pi_B$  since  $D_{KL} \geq D_{TV}^2$  according to [Schulman et al. \(2015\)](#), or an estimated maximal  $D_{TV}^2$  depending on whether the reward function is Gaussian or Categorical. We also note that for finite horizon tasks, we compute the average rewards instead of the discounted accumulated rewards in Eq. 30 and 29.

$$J_{R,1}(r; \pi_P, \pi_A) := \mathbb{E}_{\tau \sim \pi_A} [\sum_{t=0}^{\infty} \gamma^t (\delta_1(s^{(t)}, a^{(t)}) - 1) \cdot r(s^{(t)}, a^{(t)})] + C_1 \max_{(s,a) \sim \pi_A} |r(s, a)| \quad (29)$$

$$J_{R,2}(r; \pi_P, \pi_A) := \mathbb{E}_{\tau \sim \pi_P} [\sum_{t=0}^{\infty} \gamma^t (1 - \delta_2(s^{(t)}, a^{(t)})) \cdot r(s^{(t)}, a^{(t)})] + C_2 \max_{(s,a) \sim \pi_P} |r(s, a)| \quad (30)$$

Beside  $J_{R,1}$ ,  $J_{R,2}$ , we additionally use two more objective functions based on the derived bounds.  $W_{J_{R,r}}(r; \pi_A, \pi_P)$ . By denoting the optimal policy under  $r$  as  $\pi^*$ ,  $\alpha^* = \max_{s \in \mathbb{S}} D_{TV}(\pi^*(\cdot|s), \pi_A(\cdot|s))$ ,  $\epsilon^* = \max_{(s,a^{(t)})} |\mathcal{A}_{\pi^*}(s, a^{(t)})|$ , and  $\Delta \mathcal{A}_A^*(s) = \mathbb{E}_{a \sim \pi_A} [\mathcal{A}_{\pi^*}(s, a)] - \mathbb{E}_{a \sim \pi^*} [\mathcal{A}_{\pi^*}(s, a)]$ , we have the following.

$$\begin{aligned} & \eta_r(\pi_P) - \eta_r(\pi^*) \\ &= \eta_r(\pi_P) - \eta_r(\pi_A) + \eta_r(\pi_A) - \eta_r(\pi^*) \\ &\leq \eta_r(\pi_P) - \eta_r(\pi_A) + \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_A} [\Delta \mathcal{A}_A^*(s^{(t)})] + \frac{2\alpha^*\gamma\epsilon^*}{(1-\gamma)^2} \\ &= \eta_r(\pi_P) - \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_A} \left[ \mathbb{E}_{a^{(t)} \sim \pi_A} [r(s^{(t)}, a^{(t)})] \right] + \\ & \quad \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_A} \left[ \mathbb{E}_{a^{(t)} \sim \pi_A} [\mathcal{A}_{\pi^*}(s^{(t)}, a^{(t)})] - \mathbb{E}_{a^{(t)} \sim \pi^*} [\mathcal{A}_{\pi^*}(s^{(t)}, a^{(t)})] \right] + \frac{2\alpha^*\gamma\epsilon^*}{(1-\gamma)^2} \\ &= \eta_r(\pi_P) - \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_A} \left[ \mathbb{E}_{a^{(t)} \sim \pi^*} [\mathcal{A}_{\pi^*}(s^{(t)}, a^{(t)})] \right] + \frac{2\alpha^*\gamma\epsilon^*}{(1-\gamma)^2} \\ &= \mathbb{E}_{\tau \sim \pi_P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s^{(t)}, a^{(t)}) \right] - \mathbb{E}_{\tau \sim \pi_A} \left[ \sum_{t=0}^{\infty} \gamma^t \frac{\exp(r(s^{(t)}, a^{(t)}))}{\pi_A(a^{(t)}|s^{(t)})} r(s^{(t)}, a^{(t)}) \right] + \frac{2\alpha^*\gamma\epsilon^*}{(1-\gamma)^2} \quad (31) \end{aligned}$$

Let  $\delta_3 = \frac{\exp(r(s^{(t)}, a^{(t)}))}{\pi_A(a^{(t)}|s^{(t)})}$  be the importance sampling probability ratio. It is suggested in [Schulman et al. \(2017\)](#) that instead of directly optimizing the objective function Eq. 31, optimizing a surrogate

objective function as in Eq 32, which is an upper-bound of Eq 31 with some small  $\delta \in (0, 1)$  can be much less expensive and still effective.

$$J_{R,3}(r; \pi_P, \pi_A) := \mathbb{E}_{\tau \sim \pi_P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s^{(t)}, a^{(t)}) \right] - \mathbb{E}_{\tau \sim \pi_A} \left[ \sum_{t=0}^{\infty} \gamma^t \min(\delta_3 \cdot r(s^{(t)}, a^{(t)}), \text{clip}(\delta_3, 1 - \delta, 1 + \delta) \cdot r(s^{(t)}, a^{(t)})) \right] \quad (32)$$

Alternatively, we let  $\Delta \mathcal{A}_P^*(s) = \mathbb{E}_{a \sim \pi_P} [\mathcal{A}_{\pi^*}(s, a)] - \mathbb{E}_{a \sim \pi^*} [\mathcal{A}_{\pi^*}(s, a)]$ . The according to Eq 27, we have the following.

$$\begin{aligned} & \eta_r(\pi_P) - \eta_r(\pi^*) \\ & \leq \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_P} \left[ \Delta \mathcal{A}_P^*(s^{(t)}) \right] + \frac{2\alpha^* \gamma (2\alpha^* + 1) \epsilon^*}{(1 - \gamma)^2} \\ & = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s^{(t)} \sim \pi_P} \left[ \mathbb{E}_{a^{(t)} \sim \pi_P} [\mathcal{A}_{\pi^*}(s^{(t)}, a^{(t)})] - \mathbb{E}_{a^{(t)} \sim \pi^*} [\mathcal{A}_{\pi^*}(s^{(t)}, a^{(t)})] \right] + \frac{2\alpha^* \gamma (2\alpha^* + 1) \epsilon^*}{(1 - \gamma)^2} \end{aligned} \quad (33)$$

Then a new objective function  $J_{R,4}$  is formulated in Eq 34 where  $\delta_4 = \frac{\exp(r(s^{(t)}, a^{(t)}))}{\pi_P(a^{(t)} | s^{(t)})}$ .

$$J_{R,4}(r; \pi_P, \pi_A) := \mathbb{E}_{\tau \sim \pi_P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s^{(t)}, a^{(t)}) \right] - \mathbb{E}_{\tau \sim \pi_P} \left[ \sum_{t=0}^{\infty} \gamma^t \min(\delta_4 \cdot r(s^{(t)}, a^{(t)}), \text{clip}(\delta_4, 1 - \delta, 1 + \delta) \cdot r(s^{(t)}, a^{(t)})) \right] \quad (34)$$

### B.3 INCORPORATING IRL ALGORITHMS

In our implementation, we combine PAGAR with GAIL and VAIL, respectively. When PAGAR is combined with GAIL, the meta-algorithm Algorithm 1 becomes Algorithm 2. When PAGAR is combined with VAIL, it becomes Algorithm 3. Both of the two algorithms are GAN-based IRL, indicating that both algorithms use Eq 1 as the IRL objective function. In our implementation, we use a neural network to approximate  $D$ , the discriminator in Eq 1. To get the reward function  $r$ , we follow Fu et al. (2018) and denote  $r(s, a) = \log \left( \frac{\pi_A(a|s)}{D(s,a)} - \pi_A(a|s) \right)$  as mentioned in Section 1. Hence, the only difference between Algorithm 2 and Algorithm 1 is in the representation of the reward function. Regarding VAIL, since it additionally learns a representation for the state-action pairs, a bottleneck constraint  $J_{IC}(D) \leq i_c$  is added where the bottleneck  $J_{IC}$  is estimated from policy roll-outs. VAIL introduces a Lagrangian parameter  $\beta$  to integrate  $J_{IC}(D) - i_c$  in the objective function. As a result its objective function becomes  $J_{IRL}(r) + \beta \cdot (J_{IC}(D) - i_c)$ . VAIL not only learns the policy and the discriminator but also optimizes  $\beta$ . In our case, we utilize the samples from both protagonist and antagonist policies to optimize  $\beta$  as in line 10, where we follow Peng et al. (2019) by using projected gradient descent with a step size  $\delta$ .

In our implementation, depending on the difficulty of the benchmarks, we choose to maintain  $\lambda$  as a constant or update  $\lambda$  with the IRL loss  $J_{IRL}(r)$ . In *HalfCheetah-v2* and all the maze navigation tasks, we update  $\lambda$  by introducing a hyperparameter  $\mu$ . As described in the maintext, we regard  $\delta$  as the target IRL loss of  $J_{IRL}(r)$ , i.e.,  $J_{IRL}(r) \leq \delta$ . Then we update  $\lambda$  by  $\lambda := \lambda \cdot \exp(\mu \cdot (J_{IRL}(r) - \delta))$  after every iteration. In other benchmarks, we keep  $\lambda$  as a constant by letting  $\mu := 0$ . Besides, we use PPO Schulman et al. (2017) to train all policies in Algorithm 2 and 3.

## C EXPERIMENT DETAILS

This section presents some details of the experiments and additional results.

### C.1 EXPERIMENTAL DETAILS

**Network Architectures.** Our algorithm involves a protagonist policy  $\pi_P$ , and an antagonist policy  $\pi_A$ . In our implementation, the two policies have the same structures. Each structure contains two



---

**Algorithm 2** GAIL w/ PAGAR

---

**Input:** Expert demonstration  $E$ , discriminator loss bound  $\delta$ , initial protagonist policy  $\pi_P$ , antagonist policy  $\pi_A$ , discriminator  $D$  (representing  $r(s, a) = \log \left( \frac{\pi_A(a|s)}{D(s, a)} - \pi_A(a|s) \right)$ ), Lagrangian parameter  $\lambda$ , iteration number  $i = 0$ , maximum iteration number  $N$

**Output:**  $\pi_P$

```
1: while iteration number  $i < N$  do
2:   Sample trajectory sets  $\mathbb{D}_A \sim \pi_A$  and  $\mathbb{D}_P \sim \pi_P$ 
3:   Estimate  $J_{RL}(\pi_A; r)$  with  $\mathbb{D}_A$ 
4:   Optimize  $\pi_A$  to maximize  $J_{RL}(\pi_A; r)$ .
5:   Estimate  $J_{RL}(\pi_P; r)$  with  $\mathbb{D}_P$ ;  $J_{PPO}(\pi_P; \pi_A, r)$  with  $\mathbb{D}_P$  and  $\mathbb{D}_A$ ;
6:   Optimize  $\pi_P$  to maximize  $J_{RL}(\pi_P; r) + J_{PPO}(\pi_P; \pi_A, r)$ .
7:   Estimate  $J_{PAGAR}(r; \pi_P, \pi_A)$  with  $\mathbb{D}_P$  and  $\mathbb{D}_A$ 
8:   Estimate  $J_{IRL}(\pi_A; r)$  with  $\mathbb{D}_A$  and  $E$  by following the IRL algorithm
9:   Optimize  $D$  to minimize  $J_{PAGAR}(r; \pi_P, \pi_A) + \lambda \cdot \max(J_{IRL}(r) + \delta, 0)$ 
10: end while
11: return  $\pi_P$ 
```

---

---

**Algorithm 3** VAIL w/ PAGAR

---

**Input:** Expert demonstration  $E$ , discriminator loss bound  $\delta$ , initial protagonist policy  $\pi_P$ , antagonist policy  $\pi_A$ , discriminator  $D$  (representing  $r(s, a) = \log \left( \frac{\pi_A(a|s)}{D(s, a)} - \pi_A(a|s) \right)$ ), Lagrangian parameter  $\lambda$  for PAGAR, iteration number  $i = 0$ , maximum iteration number  $N$ , Lagrangian parameter  $\beta$  for bottleneck constraint, bounds on the bottleneck penalty  $i_c$ , learning rate  $\delta$ .

**Output:**  $\pi_P$

```
1: while iteration number  $i < N$  do
2:   Sample trajectory sets  $\mathbb{D}_A \sim \pi_A$  and  $\mathbb{D}_P \sim \pi_P$ 
3:   Estimate  $J_{RL}(\pi_A; r)$  with  $\mathbb{D}_A$ 
4:   Optimize  $\pi_A$  to maximize  $J_{RL}(\pi_A; r)$ .
5:   Estimate  $J_{RL}(\pi_P; r)$  with  $\mathbb{D}_P$ ;  $J_{PPO}(\pi_P; \pi_A, r)$  with  $\mathbb{D}_P$  and  $\mathbb{D}_A$ ;
6:   Optimize  $\pi_P$  to maximize  $J_{RL}(\pi_P; r) + J_{PPO}(\pi_P; \pi_A, r)$ .
7:   Estimate  $J_{PAGAR}(r; \pi_P, \pi_A)$  with  $\mathbb{D}_P$  and  $\mathbb{D}_A$ 
8:   Estimate  $J_{IRL}(\pi_A; r)$  with  $\mathbb{D}_A$  and  $E$  by following the IRL algorithm
9:   Estimate  $J_{IC}(D)$  with  $\mathbb{D}_A, \mathbb{D}_P$  and  $E$ 
10:  Optimize  $D$  to minimize  $J_{PAGAR}(r; \pi_P, \pi_A) + \lambda \cdot \max(J_{IRL}(r) + \delta, 0) + \beta \cdot J_{IC}(D)$ 
11:  Update  $\beta := \max \left( 0, \beta - \delta \cdot \left( \frac{J_{IC}(D)}{3} - i_c \right) \right)$ 
12: end while
13: return  $\pi_P$ 
```

---

neural networks, an actor network, and a critic network. When associated with GAN-based IRL, we use a discriminator  $D$  to represent the reward function as mentioned in Appendix B.3.

- **Protagonist and Antagonist policies.** We prepare two versions of actor-critic networks, an MLP version, and a CNN version, respectively, for the Mujoco and Mini-Grid benchmarks. The MLP version, the actor and critic networks have 3 layers. Each hidden layer has 100 neurons and a  $\tanh$  activation function. The output layer output the mean and standard deviation of the actions. In the CNN version, the actor and critic networks share 3 convolutional layers, each having 5, 2, 2 filters,  $2 \times 2$  kernel size, and  $ReLU$  activation function. Then 2 fully connected networks are used to simulate the actor and critic networks. The fully connected networks have one hidden layer, of which the sizes are 64.
- **Discriminator  $D$  for GAIL w/ PAGAR in Algorithm 2** We prepare two versions of discriminator networks, an MLP version and a CNN version, respectively, for the Mujoco and Mini-Grid benchmarks. The MLP version has 3 linear layers. Each hidden layer has 100 neurons and a  $\tanh$  activation function. The output layer uses the  $Sigmoid$  function to output the confidence. In the CNN version, the actor and critic networks share 3 convolutional layers, each having 5, 2, 2 filters,  $2 \times 2$  kernel size, and  $ReLU$  activation function. The last

convolutional layer is concatenated with a fully connected network with one hidden layer with 64 neurons and  $\tanh$  activation function. The output layer uses the  $Sigmoid$  function as the activation function.

- **Discriminator  $D$  for VAIL w/ PAGAR in Algorithm 3** We prepare two versions of discriminator networks, an MLP version and a CNN version, respectively, for the Mujoco and Mini-Grid benchmarks. The MLP version uses 3 linear layers to generate the mean and standard deviation of the embedding of the input. Then a two-layer fully connected network takes a sampled embedding vector as input and outputs the confidence. The hidden layer in this fully connected network has 100 neurons and a  $\tanh$  activation function. The output layer uses the  $Sigmoid$  function to output the confidence. In the CNN version, the actor and critic networks share 3 convolutional layers, each having 5, 2, 2 filters,  $2 \times 2$  kernel size, and  $ReLU$  activation function. The last convolutional layer is concatenated with a two-layer fully connected network. The hidden layer has 64 neurons and uses  $\tanh$  as the activation function. The output layer uses the  $Sigmoid$  function as the activation function.

**Hyperparameters** The hyperparameters that appear in Algorithm 3 and 3 are summarized in Table 2 where we use N/A to indicate using the maximal  $\delta$  as mentioned in Section 4.2 in which case we let  $\mu = 0$ . The values vary depending on the task and IRL algorithm for the hyperparameters  $\delta, \mu$  explained in Appendix B.3.

Parameter	Continuous Control Domain	Partially Observable Domain
Policy training batch size	64	256
Discount factor	0.99	0.99
GAE parameter	0.95	0.95
PPO clipping parameter	0.2	0.2
$\lambda$	1e3	1e3
$\sigma$	0.2	0.2
$i_c$	0.5	0.5
$\beta$	0.0	0.0
$\mu$	VAIL(HalfCheetah): 0.5; others: 0.0	VAIL: 1.0; GAIL: 1.0
$\delta$	VAIL(HalfCheetah): 1.0; others: N/A	VAIL: 0.8; GAIL: 1.2

Table 2: Hyperparameters used in the training processes

**Expert Demonstrations.** Our expert demonstrations all achieve high rewards in the task. The number of trajectories and the average trajectory total rewards are listed in Table 3.

## C.2 ADDITIONAL RESULTS

We append the results in three Mujoco benchmarks: *Hopper-v2*, *InvertedPendulum-v2* and *Swimmer-v2* in Figure 6. Algorithm 1 performs similarly to VAIL and GAIL in those two benchmarks. IQ-learn does not perform well in Walker2d-v2 but performs better than ours and other baselines by a large margin.

Task	Number of Trajectories	Average Tot.Rewards
Walker2d-v2	10	4133
HalfCheetah-v2	100	1798
Hopper-v2	100	3586
InvertedPendulum-v2	10	1000
Swimmer-v2	10	122
DoorKey-6x6-v0	10	0.92
SimpleCrossingS9N1-v0	10	0.93

Table 3: The number of demonstrated trajectories and the average trajectory rewards

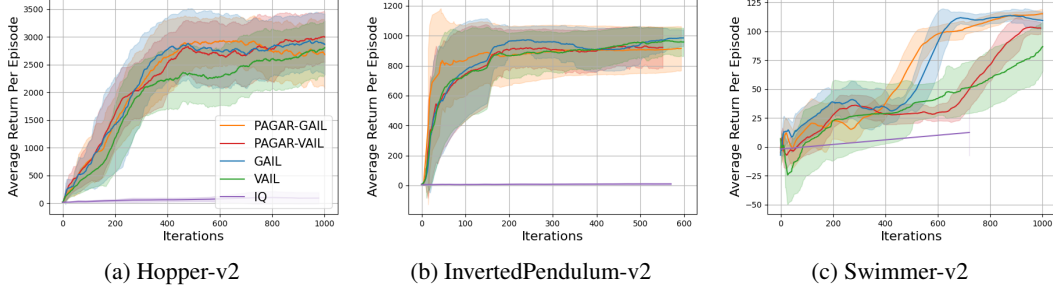


Figure 6: Comparing Algorithm 1 with baselines. The suffix after each ‘PAGAR-’ indicates which IRL algorithm is utilized in Algorithm 1. The  $y$  axis is the average return per step. The  $x$  axis is the number of iterations in GAIL, VAIL, and ours. The policy is executed between each iteration for 2048 timesteps for sample collection. One exception is that IQ-learn updates the policy at every timestep, making its actual number of iterations 2048 times larger than indicated in the figures.

### C.3 INFLUENCE OF REWARD HYPOTHESIS SPACE

We study whether choosing a different reward function set  $R$  can influence the performance of Algorithm 1. Specifically, we use the Mini-Grid benchmarks to show the difference between using the *Sigmoid* function and a Categorical distribution in the output layer of the discriminator network. When using the *Sigmoid* function, the confidence of  $D$  classifying a state-action pair as being sampled from the antagonist policy’s roll-outs is not normalized over the action space, i.e.,  $\sum_{a \in \mathbb{A}} D(s, a) \neq 1$ . When using a Categorical distribution, the confidence sum to one for all the actions, i.e.,  $\sum_{a \in \mathbb{A}} D(s, a) = 1$ . We fixed the IRL algorithm as the one used in GAIL. The results are shown in Figure 7. When changing the discriminator’s output layer, GAIL, and ours’ training efficiency change. However, our algorithm outperforms GAIL in both cases by using fewer samples to train the protagonist policy to attain high performance.

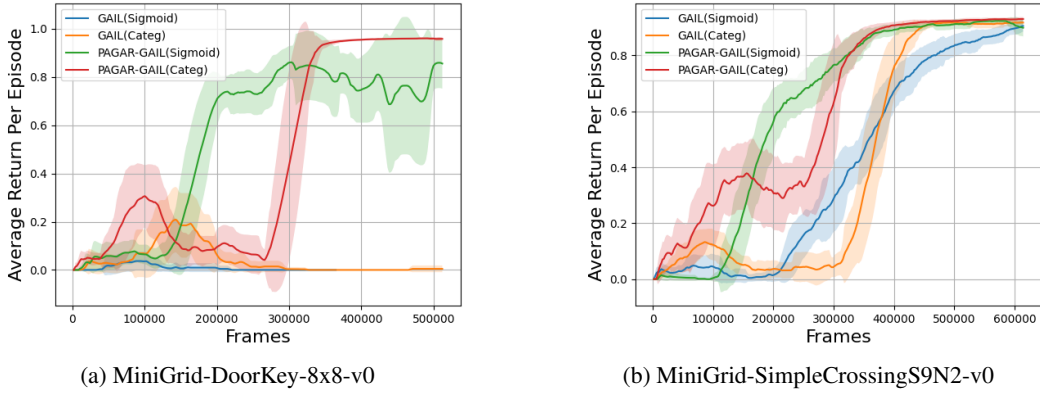


Figure 7: Comparing Algorithm 1 with baselines. The prefix ‘protagonist\_GAIL’ indicates that the IRL algorithm utilized in Algorithm 1 is the same as in GAIL. The ‘\_Sigmoid’ and ‘\_Categ’ suffixes indicate whether the output layer of the discriminator is using the *Sigmoid* function or Categorical distribution. The  $x$  axis is the number of sampled frames. The  $y$  axis is the average return per episode.