

Active Inverse Learning in Stackelberg Trajectory Games

Yue Yu, Jacob Levy, Negar Mehr, David Fridovich-Keil, and Ufuk Topcu

Abstract—Game-theoretic inverse learning is the problem of inferring the players’ objectives from their actions. We formulate an inverse learning problem in a Stackelberg game between a leader and a follower, where each player’s action is the trajectory of a dynamical system. We propose an active inverse learning method for the leader to infer which hypothesis among a finite set of candidates describes the follower’s objective function. Instead of using passively observed trajectories like existing methods, the proposed method actively maximizes the differences in the follower’s trajectories under different hypotheses to accelerate the leader’s inference. We demonstrate the proposed method in a receding-horizon repeated trajectory game. Compared with uniformly random inputs, the leader inputs provided by the proposed method accelerate the convergence of the probability of different hypotheses conditioned on the follower’s trajectory by orders of magnitude.

I. INTRODUCTION

Learning to predict human behavior is a critical challenge in human-robot interaction. It enables robots to customize their strategies in various applications, including assisted driving [1], [2], traffic management [3], [4], and, in general, mitigating conflicts in human-in-the-loop robotic systems.

Game-theoretic inverse learning helps robots explain and predict human behavior in noncooperative interactions [5], [6], [7], [8], [9], [10], [11], [3]. The idea is to first model humans’ objectives as parameterized functions, then infer the parameter value such that the corresponding game-theoretic strategies—such as Nash or Stackelberg equilibrium strategies—match the humans’ actions in a dataset. Game-theoretic inverse learning is a necessary step in understanding human-robot interactions [12], [13], [3] and designing incentives for multiagent systems [14], [15].

The existing game-theoretic inverse learning methods are *passive*, which can be data inefficient. In particular, these methods record the dataset of human actions before and independently of the inference process. Hence some actions in the recorded dataset can be uninformative for inference purposes, or simply redundant. As a result, passive inverse learning lacks the data efficiency to enable rapid inference and support online real-time decision-making.

In contrast to passive inverse learning, active inverse learning helps robots to infer human intentions in cooperative interactions by actively provoking informative human actions. For example, when learning objectives that explain human’s ranking or rating of presented options, active inverse learning

methods first provoke informative human actions and record them in the dataset, then infer the human’s objective function, and repeat this process if necessary [16], [17], [18], [19], [20]. These methods ensure that the human’s actions are informative by maximizing the volume removed from the hypothesis space [18], [21], [22], [23], [24], [25] or by maximizing the information gain [26], [27], [4], [28], [29], [30]. By integrating dataset updates with inference, active inverse learning provides practical solutions for inferring human intentions from limited interactions.

Despite its successes, active inverse learning still requires investigation in noncooperative interactions. The existing active inverse learning methods rely on querying humans who volunteer informative responses. In contrast, humans in noncooperative interactions only take actions that optimize their own objectives. Therefore, how to provoke informative actions from noncooperative humans that reveal their objectives is, to our best knowledge, still an open question.

We formulate an inverse learning problem in a Stackelberg game where a rational leader, such as a robot, is inferring which hypothesis among finitely many candidates best explains the behavior of a boundedly rational follower, such as a human. This problem is particularly relevant in shared autonomy, *e.g.*, when an autopilot is inferring the type of a newly encountered human driver. We model each player’s action as the trajectory of a linear time-invariant system. The follower tracks a linear function of the leader’s trajectory—similar to how a human driver tracks the trajectory recommended by an autopilot—using a maximum-entropy linear quadratic regulator—which contains a parameterized objective function—that models bounded rationality in human decision-making [3]. The leader determines which hypothesis is most likely using the probability of each hypothesis conditioned on the follower’s state trajectory.

We propose an active inverse learning method to provoke informative trajectories from the follower by optimizing the leader’s inputs. In this optimization, we maximize the differences in the follower’s trajectory distributions under different hypotheses. We show that this optimization is a difference-of-convex program [31], which can be solved efficiently via the convex-concave procedure [32]. We evaluate the performance of the proposed method in a receding-horizon repeated trajectory game. Compared with random inputs, the leader inputs provided by our method accelerate the convergence of the probability of different hypotheses conditioned on the follower’s trajectory by orders of magnitude.

Notation: We let \mathbb{R} , $\mathbb{R}_{\geq 0}$, and \mathbb{N} denote the set of real, non-negative real numbers, and nonnegative integers, respectively. We let $\mathbb{S}_{\geq 0}^n$ and $\mathbb{S}_{> 0}^n$ denote the set of n by n symmetric

Y. Yu, J. Levy, D. Fridovich-Keil, and U. Topcu are with the Oden Institute for Computational Engineering and Sciences, The University of Texas at Austin, Austin, TX, 78712, USA (emails: yueyu@utexas.edu, jake.levy@utexas.edu, dfk@utexas.edu, utopcu@utexas.edu). N. Mehr is with the Department of Aerospace Engineering at the University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA (email: negar@illinois.edu).

positive semidefinite and positive definite matrices, respectively. For any $x \in \mathbb{R}^n$, we let $\|x\| := \sqrt{x^\top x}$, $\|x\|_1 := \sum_{i=1}^n |x_i|$, $\|x\|_\infty := \max_{1 \leq i \leq n} |x_i|$, and $\|x\|_A^2 := x^\top A x$ for all $A \in \mathbb{S}_{>0}^n$. We let 0_n denote the n -dimensional zero vector; I_n and $0_{n \times n}$ denote the n by n identity and zero matrix, respectively. We let $\mathcal{N}(\mu, \Sigma)$ denote the Gaussian distribution with mean $\mu \in \mathbb{R}^n$ and variance $\Sigma \in \mathbb{S}_{>0}^n$. Given $n_1, n_2 \in \mathbb{N}$, we let $[n_1, n_2]$ denote the set of integers between n_1 and n_2 . Given $a_i \in \mathbb{R}^n$ for all $i \in \mathbb{N}$, we let $a_{i:j} := [a_i^\top \ a_{i+1}^\top \ \dots \ a_j^\top]^\top$ for all $i < j$, $i, j \in \mathbb{N}$.

II. LINEAR QUADRATIC STACKELBERG TRAJECTORY GAMES

We introduce a Stackelberg game between a rational leader, such as a robot, and a boundedly rational follower, such as a human with noisy behavior. The players' actions are trajectories of stochastic linear time-invariant systems.

A. The dynamics of the players' systems

We assume that the leader's state evolves according to the following discrete-time linear time-invariant dynamics:

$$x_{t+1}^L = A^L x_t^L + B^L u_t^L + w_t^L \quad (1)$$

for all $t \in \mathbb{N}$, where $x_t^L \in \mathbb{R}^{n_L}$, $u_t^L \in \mathbb{R}^{m_L}$, and $w_t^L \in \mathbb{R}^{n_L}$ are the state, input, and disturbance of the system at time $t \in \mathbb{N}$, respectively; $A^L \in \mathbb{R}^{n_L \times n_L}$ and $B^L \in \mathbb{R}^{n_L \times m_L}$ are the leader's system parameters. Equation (1) characterizes the dynamics of many common robots, such as the kinematics of rovers and drones.

Similarly, the follower's state evolves according to the following dynamics:

$$x_{t+1}^F = A^F x_t^F + B^F u_t^F + w_t^F \quad (2)$$

for all $t \in \mathbb{N}$, where $x_t^F \in \mathbb{R}^{n_F}$, $u_t^F \in \mathbb{R}^{m_F}$, and $w_t^F \in \mathbb{R}^{n_F}$ denote the state, input, and disturbance of the system at time $t \in \mathbb{N}$, respectively; $A^F \in \mathbb{R}^{n_F \times n_F}$ and $B^F \in \mathbb{R}^{n_F \times m_F}$ are the follower's system parameters.

Throughout, we assume that the disturbance in the leader and the follower's systems are independent, identically distributed Gaussian vectors, *i.e.*, there exists $\Omega^L \in \mathbb{S}_{>0}^{n_L}$ and $\Omega^F \in \mathbb{S}_{>0}^{n_F}$ such that, for any $t \in \mathbb{N}$, we have

$$w_t^L \sim \mathcal{N}(0_{n_L}, \Omega^L), \quad w_t^F \sim \mathcal{N}(0_{n_F}, \Omega^F). \quad (3)$$

B. The players' objectives

We assume that the follower's objective is to track a linear function of the leader's trajectory. In particular, we let $M^F \in \mathbb{R}^{n_F \times n_L}$ denote a matrix that maps the leader's internal state to an output reference observable to the follower. Let $x_{0:\tau}^L$ denote a leader trajectory of length $\tau \in \mathbb{N}$. The follower's objective is to simultaneously track the corresponding output trajectory $\{M^F x_0^L, M^F x_1^L, \dots, M^F x_\tau^L\}$ and minimize its input efforts.

We assume that the follower is boundedly rational and chooses its input according to the maximum entropy principle, which states that the distribution of u_t^F conditioned on x_t^F is Gaussian, *i.e.*, $\mu_t^F | x_t^F \sim \mathcal{N}(\mu_t, \Sigma_t)$ for some $\mu_t \in \mathbb{R}^{m_F}$

and $\Sigma_t \in \mathbb{R}^{m_F}$ [3]. In particular, $(\mu_{0:\tau-1}, \Sigma_{0:\tau-1})$ is optimal for the following stochastic trajectory optimization problem:

$$\begin{aligned} & \underset{\mu_{0:\tau-1}, \Sigma_{0:\tau-1}}{\text{minimize}} && \sum_{t=0}^{\tau} \mathbb{E} \left[\frac{1}{2} \|x_t^F - M^F x_t^L\|_{Q^F}^2 \right] \\ & \text{subject to} && + \frac{1}{2} \sum_{t=0}^{\tau-1} \left(\mathbb{E} \left[\|u_t^F\|_{R^F}^2 \right] - \log \det \Sigma_t \right) \\ & && x_{t+1}^F = A^F x_t^F + B^F u_t^F + w_t^F, \quad x_0^F = \hat{x}_0^F, \\ & && u_t^F | x_t^F \sim \mathcal{N}(\mu_t, \Sigma_t), \quad w_t^F \sim \mathcal{N}(0_{n_F}, \Omega^F), \\ & && t \in [0, \tau - 1], \end{aligned} \quad (4)$$

where $\hat{x}_0^F \in \mathbb{R}^F$ is the initial state of the follower's system, $\mathbb{E}[\cdot]$ denotes the expectation; $Q^F \in \mathbb{S}_{>0}^{n_F}$, and $R^F \in \mathbb{S}_{>0}^{m_F}$ are the follower's cost parameters. The objective function in optimization (4) captures boundedly rational human decisions: it is noisy but centers around a cost-minimizing rational decision.

The following proposition provides a closed-form formula for the solution of optimization (4).

Proposition 1. Let

$$P_\tau^F = Q^F, \quad P^F = Q^F + (A^F)^\top P_{t+1}^F E_t^F, \quad (5a)$$

$$F_t^F = B^F (R + (B^F)^\top P_{t+1}^F B^F)^{-1} (B^F)^\top, \quad (5b)$$

$$E_t^F = A^F - F_t^F P_{t+1}^F A^F, \quad (5c)$$

$$q_\tau^F = -Q^F M^F x_\tau^L, \quad q_t^F = (E_t^F)^\top q_{t+1}^F - Q^F M^F x_t^L, \quad (5d)$$

for all $t \in [0, \tau - 1]$. Given $x_{0:\tau}^L$, $(\mu_{0:\tau-1}, \Sigma_{0:\tau-1})$ is optimal for optimization (4) if and only if

$$\Sigma_t = (R^F + (B^F)^\top P_{t+1}^F B^F)^{-1}, \quad (6a)$$

$$\mu_t = -\Sigma_t (B^F)^\top (P_{t+1}^F A^F x_t^F + q_{t+1}^F). \quad (6b)$$

for all $t \in [0, \tau - 1]$. Furthermore, if the constraints in (4) hold, then $x_t^F \sim \mathcal{N}(\xi_t, \Lambda_t)$ for all $t \in [0, \tau]$, where

$$\xi_{t+1} = E_t^F \xi_t - F_t^F q_{t+1}^F, \quad (7a)$$

$$\Lambda_{t+1} = E_t^F \Lambda_t (E_t^F)^\top + F_t^F + \Omega^F, \quad (7b)$$

for all $t \in [0, \tau - 1]$, with $\xi_0 = \hat{x}_0^F$ and $\Lambda_0 = 0_{n_F \times n_F}$.

Proof. See Appendix. \square

The leader's objective is to minimize a cost function that jointly depends on the expected follower's trajectory and the leader's trajectory. To this end, we assume that the leader acts rationally and chooses its trajectory $x_{0:\tau}^L$ as a solution to the following trajectory optimization problem:

$$\begin{aligned} & \underset{u_{0:\tau-1}}{\text{minimize}} && \mathbb{E}[f(x_{0:\tau}^L, x_{0:\tau}^F)] + g(u_{0:\tau}^L) \\ & \text{subject to} && x_{t+1}^L = A^L x_t^L + B^L u_t^L + w_t^L, \quad x_0^L = \hat{x}_0^L, \\ & && x_{t+1}^F = A^F x_t^F + B^F u_t^F + w_t^F, \quad x_0^F = \hat{x}_0^F, \\ & && w_t^L \sim \mathcal{N}(0_{n_L}, \Omega^L), \quad w_t^F \sim \mathcal{N}(0_{n_F}, \Omega^F), \\ & && u_t \in \mathbb{U}, \quad u_t^F | x_t^F \sim \mathcal{N}(\mu_t, \Sigma_t), \quad t \in [0, \tau - 1], \\ & && (\mu_{0:\tau-1}, \Sigma_{0:\tau-1}) \text{ is optimal for (4),} \end{aligned} \quad (8)$$

where $\hat{x}_0^L \in \mathbb{R}^{n_L}$ is the initial state of the leader's system, $\mathbb{U} \subset \mathbb{R}^{m_F}$ is the set of feasible leader inputs at each time. Furthermore, $f : \mathbb{R}^{(\tau+1)n_L} \times \mathbb{R}^{(\tau+1)n_F} \rightarrow \mathbb{R}$ is the leader's cost function that jointly depends on the leader's state trajectory $x_{0:\tau}^L$ and the follower's state trajectory $x_{0:\tau}^F$;

$g : \mathbb{R}^{\tau m_L} \rightarrow \mathbb{R}$ is a cost function that only depends on the leader's input trajectory. By choosing different functions for f and g , optimization (8) achieves different trade-offs between optimizing the leader and the follower's trajectory.

Problem (8) is a Stackelberg game, also known as a *bilevel optimization*. See [33] and reference therein for a detailed discussion on bilevel optimization.

III. ACTIVE INVERSE LEARNING VIA DIFFERENCE MAXIMIZATION

Given the Stackelberg trajectory game introduced in Section II, we now consider the case where the leader does not know the parameter tuple (Q^F, R^F, M^F) in the follower's objective, except that it is one of finitely many candidates. In other words, the leader knows that there exist $Q^1, \dots, Q^d \in \mathbb{R}^{n_F \times n_F}$, $R^1, \dots, R^d \in \mathbb{R}^{m_F \times m_F}$, and $M^1, \dots, M^d \in \mathbb{R}^{n_F \times n_L}$ such that

$$(Q^F, R^F, M^F) = (Q^i, R^i, M^i) \quad (9)$$

for some $i \in [1, d]$. This case arises, for example, when a robot already learned different types of human behavior offline but needs to determine the type of a newly encountered human via online interaction. In the following, we let $\theta^F := (Q^F, R^F, M^F)$ and $\theta^i := (Q^i, R^i, M^i)$ for all $i \in [1, d]$. We say that *hypothesis i is true* if (9) holds.

Based on a prior probability distribution of all hypotheses that gives the value of $\mathbb{P}(\theta^F = \theta^i | x_0^F)$ for all $i \in [1, d]$ and the follower's trajectory $x_{1:\tau}^F$, the leader can infer whether hypothesis i is more likely than hypothesis j to be true by computing the following ratio:

$$\frac{\mathbb{P}(\theta^F = \theta^i | x_{0:\tau}^F)}{\mathbb{P}(\theta^F = \theta^j | x_{0:\tau}^F)} = \frac{\mathbb{P}(\theta^F = \theta^i | x_0^F) \mathbb{P}(x_{1:\tau}^F | \theta^F = \theta^i)}{\mathbb{P}(\theta^F = \theta^j | x_0^F) \mathbb{P}(x_{1:\tau}^F | \theta^F = \theta^j)} \quad (10)$$

The ratio in (10) being bigger than one implies that trajectory $x_{0:\tau}^F$ is more likely to occur under hypothesis i rather than hypothesis j , and vice versa.

However, observing the follower's trajectory can be uninformative for the inference above if the trajectories under different hypotheses are similar. For example, suppose that

$$\mathbb{P}(x_{1:\tau}^F | \theta^F = \theta^i) \approx \mathbb{P}(x_{1:\tau}^F | \theta^F = \theta^j) \quad (11)$$

for some $i \neq j$, then (10) implies that $\frac{\mathbb{P}(\theta^F = \theta^i | x_{0:\tau}^F)}{\mathbb{P}(\theta^F = \theta^j | x_{0:\tau}^F)} \approx \frac{\mathbb{P}(\theta^F = \theta^i | x_0^F)}{\mathbb{P}(\theta^F = \theta^j | x_0^F)}$. In other words, observing the follower's ongoing trajectory $x_{1:\tau}^F$ does not help the leader distinguish hypothesis i from j . In the following, we discuss how the leader can actively avoid the scenarios where (11) happens by making the follower's trajectories under different hypotheses as different as possible.

A. The follower's trajectories under different hypotheses

To avoid the scenarios where (11) happens, it suffices to make the follower's trajectory distributions under different hypotheses as different as possible. To this end, we first take a closer look at the follower's trajectory. Proposition 1 shows that the follower's state at each time is a Gaussian random variable. Particularly, let

$$E_t^i = A^F - F_t^i P_{t+1}^i A^F, \quad (12a)$$

$$F_t^i = B^F (R^F + (B^F)^\top P_{t+1}^i B^F)^{-1} (B^F)^\top, \quad (12b)$$

$$P_\tau^i = Q^i, P_t^i = Q^i + (A^F)^\top P_{t+1}^i E_t^i, \quad (12c)$$

$$\Lambda_0^i = 0_{n_F \times n_F}, \Lambda_{t+1}^i = E_t^i \Lambda_t^i (E_t^i)^\top + F_t^i + \Omega^F, \quad (12d)$$

$$q_\tau^i = -Q^i M^i x_\tau^L, q_t^i = (E_t^i)^\top q_{t+1}^i - Q^i M^i x_t^L, \quad (12e)$$

$$\xi_0^i = \hat{x}_0^F, \xi_{t+1}^i = E_t^i \xi_t^i + F_t^i q_{t+1}^i, \quad (12f)$$

for all $t \in [0, \tau - 1]$. Proposition 1 implies that, if hypothesis i is true, i.e., $\theta^F = \theta^i$, then

$$x_t^F \sim \mathcal{G}_t^i := \mathcal{N}(\xi_t^i, \Lambda_t^i). \quad (13)$$

To measure the differences between the trajectory distribution under different hypotheses, we introduce a distance function. To this end, let

$$\mathbb{D} := \{(i, j) | i < j, i, j \in [1, d]\}. \quad (14)$$

A popular measure of the difference between two Gaussian distributions is the *KL-divergence*. Given $(i, j) \in \mathbb{D}$ and $t \in [1, \tau]$, let \mathcal{G}_t^i and \mathcal{G}_t^j be defined as in (13). The KL-divergence from \mathcal{G}_t^i to \mathcal{G}_t^j is as follows:

$$D_{KL}(\mathcal{G}_t^i || \mathcal{G}_t^j) := \frac{1}{2} \left\| \xi_t^i - \xi_t^j \right\|_{(\Lambda_t^i)^{-1}}^2 - \frac{(\tau+1)n_F}{2} + \frac{1}{2} \log \left(\frac{\det \Lambda_t^j}{\det \Lambda_t^i} \right) + \text{tr}((\Lambda_t^j)^{-1} \Lambda_t^i). \quad (15)$$

Notice that, since $\Omega^F \in \mathbb{S}_{>0}^{n_F}$, (12) implies that $\Lambda_t^i \in \mathbb{S}_{>0}^{n_F}$ for all $t \in [1, \tau]$ and $k \in [1, d]$. However, KL divergence is not symmetric, i.e., $D_{KL}(\mathcal{G}_t^i || \mathcal{G}_t^j) \neq D_{KL}(\mathcal{G}_t^j || \mathcal{G}_t^i)$. To define a symmetric distance function, we propose the following function

$$D(\mathcal{G}_t^i, \mathcal{G}_t^j) := \left\| \xi_t^i - \xi_t^j \right\|_{(\Lambda_t^i)^{-1} + (\Lambda_t^j)^{-1}}^2 \quad (16)$$

for all $t \in [1, \tau]$ and $(i, j) \in \mathbb{D}$. The intuition behind this distance function is to first evaluate (up to a constant of 2, for the convenience of notation) the sum of $D_{KL}(\mathcal{G}_t^i || \mathcal{G}_t^j)$ and $D_{KL}(\mathcal{G}_t^j || \mathcal{G}_t^i)$, then remove the terms that are independent of ξ_t^i and ξ_t^j , which are, as suggested by (12), independent of the leader's trajectory. Later, we use this function to optimize the leader's inputs.

B. Maximizing the worst-case pairwise distance

To avoid the scenarios where (11) happens, we need to maximize the value of the distance function in (16) for any $(i, j) \in \mathbb{D}$. To this end, define the following *worst-case distance function*, which evaluates the minimum value of function (16) among all $(i, j) \in \mathbb{D}$:

$$\begin{aligned} & \min_{(i,j) \in \mathbb{D}} \left\{ \sum_{t=1}^{\tau} \left\| \xi_t^i - \xi_t^j \right\|_{(\Lambda_t^i)^{-1} + (\Lambda_t^j)^{-1}}^2 \right\} \\ &= \sum_{(i,j) \in \mathbb{D}} \sum_{t=1}^{\tau} \left\| \xi_t^i - \xi_t^j \right\|_{(\Lambda_t^i)^{-1} + (\Lambda_t^j)^{-1}}^2 \\ &= \max_{(i,j) \in \mathbb{D}} \left\{ \sum_{(k,l) \in \mathbb{D} \setminus \{(i,j)\}} \sum_{t=1}^{\tau} \left\| \xi_t^k - \xi_t^l \right\|_{(\Lambda_t^k)^{-1} + (\Lambda_t^l)^{-1}}^2 \right\}. \end{aligned} \quad (17)$$

The second step in (17) uses the fact that, given any $\alpha_1, \dots, \alpha_n \in \mathbb{R}$, we have

$$\min_{i \in [1, n]} \alpha_i = \sum_{i \in [1, n]} \alpha_i - \max_{i \in [1, n]} \sum_{j \in [1, n], j \neq i} \alpha_j.$$

Notice that the distance in (17) does not include the terms for $t = 0$ since, due to (12), $\xi_0^i = \xi_0^j$ for all $(i, j) \in \mathbb{D}$.

Based on the worst-case distance in (17), we propose to optimize the leader's input trajectory $u_{0:\tau-1}^L$ via solving optimization (18) (see next page) instead of optimization (8), where E_t^i, F_t^i and Λ_t^i are given by (12). The idea of optimization (18) is to first approximate the leader's state trajectory $x_{0:\tau}^L$ with its expectation, denoted by $\eta_{0:\tau}$, which satisfies the averaged dynamics $\eta_{t+1} = A^L \eta_t + B^L u_t^L$. This approximation replaces the disturbance term w_t^L in (8) with its most probable estimation of its distribution $\mathcal{N}(0_{n_L}, \Omega^L)$, given by $w_t^L = 0_{n_L}$. In addition, optimization (18) replaces the $\mathbb{E}[f(x_{0:\tau}^L, x_{0:\tau}^F)]$ term in (8) with the negative of the worst-case distance in (17). In particular, the constraints in (18) imply that

$$\max_{(i,j) \in \mathbb{D}} \left\{ \sum_{(k,l) \in \mathbb{D} \setminus \{(i,j)\}} \sum_{t=1}^{\tau} \|\xi_t^k - \xi_t^l\|_{(\Lambda^k)_t^{-1} + (\Lambda^l)_t^{-1}}^2 \right\} \leq s.$$

Since the objective function in (18) is minimizing the value of s , the above inequality holds as an equality at optimality. Hence the objective function in (18) is equivalent to the one in (8) except that the $\mathbb{E}[f(x_{0:\tau}^L, x_{0:\tau}^F)]$ term in (8) is replaced by the negative of the worst-case distance in (17). The idea of this replacement is to maximize the worst-case distance in (17), which ensures that the distance in (16) is large for any $(i, j) \in \mathbb{D}$.

Optimization (18) is a difference-of-convex program: all of its constraints are convex, but its objective function is the difference between two convex functions [31]. A popular solution method for difference-of-convex programs is the *convex-concave procedure*, which guarantees global convergence to a stationary point [34], and provides locally optimal solutions in practice [32].

IV. NUMERICAL EXPERIMENTS

We empirically demonstrate our results using a receding-horizon repeated trajectory game between a boundedly rational follower that controls one ground rover and a rational leader that controls multiple ground rovers. The leader is inferring which leading rover the follower is following. At each time step, the players play a Stackelberg trajectory game and implement only the first step of their respective input trajectories.

We define the game parameters as follows. We model the dynamics of the following rover using an instance of (1), where $A^F = \exp\left(\delta \begin{bmatrix} 0_{2 \times 2} & I_2 \\ 0_{2 \times 2} & 0_{2 \times 2} \end{bmatrix}\right)$, $B^F = \int_0^\delta \exp\left(t \begin{bmatrix} 0_{2 \times 2} & I_2 \\ 0_{2 \times 2} & 0_{2 \times 2} \end{bmatrix}\right) dt \begin{bmatrix} 0_{2 \times 2} \\ I_2 \end{bmatrix}$, $\Omega^F = \frac{1}{1000} I_4$, and $\delta = 0.2$ is the discretization step size. For the leader, we model the joint dynamics of $d \in \mathbb{N}$ ground rovers using the joint dynamics of d double-integrators, where we let $A^L = I_d \otimes A$, $B^L = I_d \otimes B$, and $\Omega^L = \frac{1}{1000} I_{4d \times 4d}$. For the follower's trajectory optimization in (4), we let $Q^i = \text{diag}([1 \ 1 \ 0 \ 0])$, $R^i = \frac{1}{100} I_2$, and $M^i = [0_{4 \times 4(i-1)} \ I_4 \ 0_{4 \times 4(d-i)}]$ for all $i \in [1, d]$. Without loss of generality, we assume $(Q^F, R^F, M^F) = (Q^1, R^1, M^1)$. For the leader's trajectory optimization in (8), we let $\mathbb{U} := \{u \in \mathbb{R}^{2d} \mid \|u\|_\infty \leq 2\}$, and $g(u_{0:\tau-1}^L) = \sum_{t=1}^{\tau-2} \|u_{t+1}^L - u_t^L\|^2$.

We demonstrate the proposed learning methods using the simulation of a receding-horizon repeated trajectory game. At time $t\delta$ for some $t \in \mathbb{N}$, the leader first observes its current state \bar{x}_t^L and the follower's current state \bar{x}_t^F , then solves optimization (18) with $\hat{x}_0^L = \bar{x}_t^L$ and $\hat{x}_0^F = \bar{x}_t^F$ to obtain the optimal input sequence $u_{0:\tau-1}^L$. Next, the leader simulates a state trajectory $x_{0:\tau-1}^L$ according to (1) and shares it with the follower, then applies u_0^L . Meanwhile, at time $t\delta$, the follower observes its current state \hat{x}_t^F and receives the leader's simulated trajectory $x_{0:\tau}^L$. Next the follower solves optimization (4) with $\hat{x}_0^F = \bar{x}_t^F$ and obtain the optimal mean and covariance sequences $(\mu_{0:\tau-1}, \Sigma_{0:\tau-1})$, and finally samples $u_0^F \sim \mathcal{N}(\mu_0, \Sigma_0)$ and applies u_0^F to its system.

We simulate the players' trajectories in this receding-horizon repeated trajectory game, where we solve the leader's trajectory optimization (18) using the convex-concave procedure [32]. Fig. 1 shows the position trajectories of the leader's rovers when $d = 3$ and $d = 5$. We can see that, to make the follower's trajectories under different hypotheses as different as possible, optimization (18) ensures that different leading rovers are moving in different directions.

We further showcase the advantage of the proposed method in terms of distinguishing different hypotheses. To this end, we let $p^* = [1 \ 0 \ \dots \ 0]^T \in \mathbb{R}^d$ denote the ground truth probability of all hypotheses (note that we let $(Q^F, R^F, M^F) = (Q^1, R^1, M^1)$). In addition, we let p^t denote the d -dimensional vector such that

$$p_i^t := \mathbb{P}((Q^F, R^F, M^F) = (Q^i, R^i, M^i) \mid \bar{x}_{0:t}^F) \quad (19)$$

for all $i \in [1, d]$, where we let $\mathbb{P}((Q^F, R^F, M^F) = (Q^i, R^i, M^i) \mid \bar{x}_{0:t}^F) := \frac{1}{d}$ for all $i \in [1, d]$, i.e., we choose the uniform distribution as the leader's prior distribution over all hypotheses given the follower's initial state. Notice that one can compute p^t recursively using the Bayes rule. Fig. 2 shows the time history of the ℓ_1 -norm distance between p^t and p^* , and compare the results where the leader uses, instead of the inputs optimal for optimization (18), independent and identically distributed input sampled from the uniform distribution over \mathbb{U} . The results in Fig. 2 show that, when compared with uniformly random inputs, the proposed method provide leader inputs that reduce the difference between p^t and p^* by orders of magnitudes.

V. CONCLUSION

We formulated an inverse learning problem in a Stackelberg trajectory game, where the leader is inferring the type of the follower's cost function by observing its trajectories. We proposed an active inverse learning method to accelerate the leader inference by making the follower's trajectories under different hypotheses as different as possible. This method accelerates the convergence of the probability of different hypotheses by orders of magnitude when compared against random inputs.

However, the current work still has limitations. For example, it considers neither nonlinear dynamics nor constraints for input limits and collision avoidance in the players' trajectory optimization. In addition, it ignores the possibility

$$\begin{aligned}
& \underset{s, u_{0:\tau-1}^i}{\text{minimize}} && s - \sum_{(i,j) \in \mathbb{D}} \sum_{t=0}^{\tau} \left\| \xi_t^i - \xi_t^j \right\|_{(\Lambda_i^i)^{-1} + (\Lambda_j^j)^{-1}}^2 + g(u_{0:\tau-1}) \\
& \text{subject to} && \eta_{t+1}^L = A^L \eta_t^L + B^L u_t^L, \eta_0^L = \hat{x}_0^L, \\
& && q_t^i = (E_t^i)^\top q_{t+1}^i - Q^i M^i \eta_t, q_\tau^i = -Q^i M^i \eta_\tau, \xi_{t+1}^i = E_t^i \xi_t^i - F_t^i q_{t+1}^i, \xi_0^i = \hat{x}_0^i, \\
& && \sum_{(k,l) \in \mathbb{D} \setminus \{(i,j)\}} \sum_{t=1}^{\tau} \left\| \xi_t^k - \xi_t^l \right\|_{(\Lambda^k)^{-1} + (\Lambda^l)^{-1}}^2 \leq s, \forall t \in [0, \tau-1], i \in [1, d], (i,j) \in \mathbb{D}.
\end{aligned} \tag{18}$$

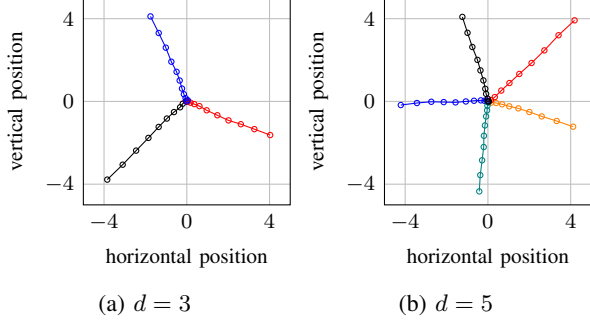


Fig. 1: The position trajectories of the leader's rovers, where different rover's trajectories are marked by different colors.

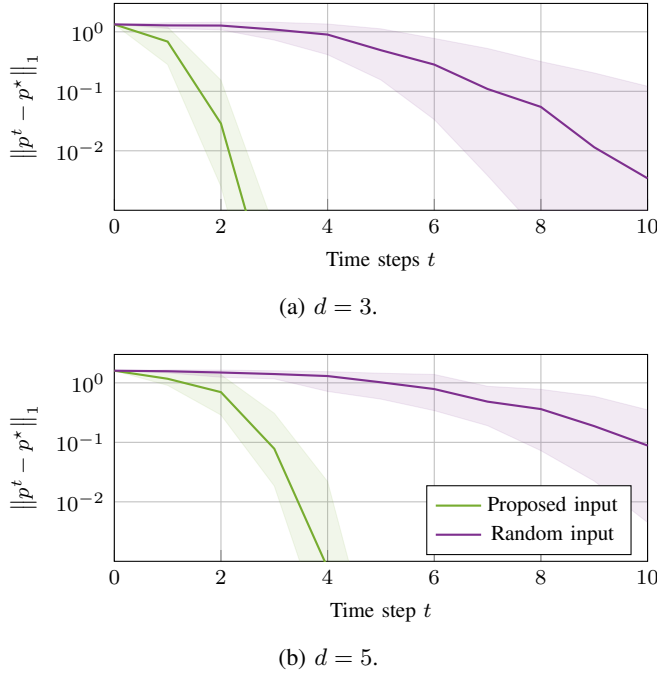


Fig. 2: The comparison of the convergence of p^t when the leader obtains its input by solving optimization 18 versus sampling uniformly in set \mathbb{U} . The solid lines show the median over 100 simulations, while the boundaries of the colored areas mark the corresponding first and third quartiles.

of deceptive actions of the follower. In future work, we plan to address these limitations and develop information-gathering strategies in general Bayesian games.

APPENDIX

Proof of Proposition 1

Given $y \in \mathbb{R}^{n_F}$ and $t \in [0, \tau]$, let

$$\begin{aligned}
& \psi_t(y, \mu_{t:\tau}, \Sigma_{t:\tau}) \\
& := \sum_{j=t}^{\tau} \mathbb{E} \left[\frac{1}{2} \|x_j^F - M^F x_j^L\|_{Q^F}^2 \mid x_t^F = y \right] \\
& + \frac{1}{2} \sum_{j=t}^{\tau-1} \left(\mathbb{E} \left[\|u_j^F\|_{R^F}^2 \mid x_t^F = y \right] - \log \det \Sigma_j \right).
\end{aligned} \tag{20}$$

where $x_{j+1}^F = A^F x_j^F + B^F u_j^F$ and $u_j^F \sim \mathcal{N}(\mu_j, \Sigma_j)$. Furthermore, let

$$V_t(y) := \min_{\mu_{t:\tau-1}, \Sigma_{t:\tau-1}} \psi_t(y, \mu_{t:\tau}, \Sigma_{t:\tau}) \tag{21}$$

for all $t \in [0, \tau]$. Then one can verify that the optimal value of optimization (4) is $V(\hat{x}_0^F)$. In addition, we can show that $V_\tau(y) = \frac{1}{2} y^\top Q^F y - \langle Q^F M^F x_\tau^L, y \rangle + \frac{1}{2} \|M^F x_\tau^L\|_{Q^F}^2$, i.e., $V_\tau(y)$ is a quadratic function of y . Suppose that $V_{t+1}(y)$ is a quadratic function of y , i.e., there exists $P_{t+1}^F \in \mathbb{R}^{n_F \times n_F}$, $q_{t+1}^F \in \mathbb{R}^{n_F}$, and $\nu_{t+1}^F \in \mathbb{R}$, such that

$$V_{t+1}(y) = \frac{1}{2} y^\top P_{t+1}^F y + \langle q_{t+1}^F, y \rangle + \nu_{t+1}^F. \tag{22}$$

Then (21) and the principle of dynamic programming together imply that

$$\begin{aligned}
& V_t(y) \\
& = \min_{\mu_t, \Sigma_t} \mathbb{E} \left[\frac{1}{2} \|u_t^F\|_{R^F}^2 + V_{t+1}(A^F x_t^F + B^F u_t^F + w_t^F) \mid x_t^F = y \right] \\
& + \frac{1}{2} \|y - M^F x_t^F\|_{Q^F}^2 - \frac{1}{2} \log \det \Sigma_t
\end{aligned} \tag{23}$$

where $u_t^F \mid x_t^F \sim \mathcal{N}(\mu_t, \Sigma_t)$. Observe that

$$\begin{aligned}
& \mathbb{E}[\|u_t^F\|_{R^F}^2 \mid x_t^F = y] \\
& = \mu_t^\top R^F \mu_t + \mathbb{E}[\text{tr}((u_t^F - \mu_t)(u_t^F - \mu_t)^\top R^F) \mid x_t^F = y] \\
& = \mu_t^\top R^F \mu_t + \text{tr}(\Sigma_t R^F).
\end{aligned} \tag{24}$$

In addition, by using (22) we can show that

$$\begin{aligned}
& \mathbb{E}[V_{t+1}(A^F x_t^F + B^F u_t^F + w_t^F) \mid x_t^F = y] \\
& = \frac{1}{2} (A^F y + B^F \mu_t)^\top P_{t+1}^F (A^F y + B^F \mu_t) \\
& + \mathbb{E}[(A^F x_t^F + B^F \mu_t)^\top P_{t+1}^F B^F (u_t^F - \mu_t) \mid x_t^F = y] \\
& + \frac{1}{2} \mathbb{E}[(B^F (u_t^F - \mu_t))^\top P_{t+1}^F B^F (u_t^F - \mu_t) \mid x_t^F = y] \\
& + \langle q_{t+1}^F, A^F y + B^F \mu_t \rangle + \frac{1}{2} \mathbb{E}[\text{tr}(w_t^F (w_t^F)^\top P_{t+1}^F)] + \nu_{t+1}^F \\
& = \frac{1}{2} (A^F y + B^F \mu_t)^\top P_{t+1}^F (A^F y + B^F \mu_t) + \frac{1}{2} \text{tr}(\Omega^F P_{t+1}^F) \\
& + \frac{1}{2} \text{tr}(\Sigma_t (B^F)^\top P_{t+1}^F B^F) + \langle q_{t+1}^F, A^F y + B^F \mu_t \rangle + \nu_{t+1}^F.
\end{aligned} \tag{25}$$

Substituting (25) and (24) into (23) gives

$$\begin{aligned}
V_t(y) &= \frac{1}{2}y^\top(Q^F + (A^F)^\top P_{t+1}^F A^F)y + \langle A^\top q_{t+1}^F - Q^F M^F x_t^L, y \rangle \\
&+ \frac{1}{2}\mu_t^\top(R^F + (B^F)^\top P_{t+1}^F B^F)\mu_t + \frac{1}{2}(x_t^L)^\top(M^F)^\top Q^F M^F x_t^L \\
&+ \langle (B^F)^\top q_{t+1}^F + (B^F)^\top P_{t+1}^F A^F y, \mu_t \rangle - \frac{1}{2} \log \det \Sigma_t \\
&+ \frac{1}{2} \text{tr}(\Sigma_t R^F + \Omega^F P_{t+1} + (B^F)^\top P_{t+1}^F B^F) + \nu_{t+1}^F.
\end{aligned} \tag{26}$$

By setting the derivative of $V_t(y)$ with respect to μ_t and Σ_t to zero, we obtain (6). By substituting (6) into (26), we can show that $V_t(y) = \frac{1}{2}y^\top P_t^F y + \langle q_t^F, y \rangle + \nu_t^F$, where Q_t^F and q_t^F satisfy (5).

Next, let $K_t = -\Sigma_t(B^F)^\top P_{t+1}^F A^F$ and $b_t = -\Sigma_t(B^F)^\top q_{t+1}^F$. Then (6) implies that $u_t^F | x_t^F \sim \mathcal{N}(K_t x_t^F + b_t, \Sigma_t)$. Since $x_0^F = \hat{x}_0^F$, by using the results in [35, p. 91] we can show the following:

$$\begin{bmatrix} x_1^F \\ u_1^F \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \xi_0 \\ K_0 \xi_0 + b_0 \end{bmatrix}, \begin{bmatrix} \Lambda_0 & K_0 \Lambda_0 \\ \Lambda_0 K_0^\top & \Sigma_0 + K_0^\top \Lambda_0 K_0 \end{bmatrix} \right).$$

Therefore $x_{t+1}^F = A^F x_t^F + B^F u_t^F + w_t^F \sim \mathcal{N}(\xi_1, \Lambda_1)$, where ξ_1 and Σ_1 satisfy (7). By repeating similar steps for $t \in [2, \tau]$ we can show that (7) holds for all $t \in [0, \tau - 1]$, which completes the proof.

REFERENCES

- [1] V. A. Shia, Y. Gao, R. Vasudevan, K. D. Campbell, T. Lin, F. Borrelli, and R. Bajcsy, "Semiautonomous vehicular control using driver modeling," *IEEE Trans. Intel. Transp. Syst.*, vol. 15, no. 6, pp. 2696–2709, 2014.
- [2] N. Mehr, R. Horowitz, and A. D. Dragan, "Inferring and assisting with constraints in shared autonomy," in *Proc. IEEE Conf. Decision Control*. IEEE, 2016, pp. 6689–6696.
- [3] N. Mehr, M. Wang, M. Bhatt, and M. Schwager, "Maximum-entropy multi-agent dynamic games: Forward and inverse solutions," *IEEE Trans. Robot.*, 2023.
- [4] D. Sadigh, S. S. Sastry, S. A. Seshia, and A. Dragan, "Information gathering actions over human internal state," in *Proc. Int. Conf. Intel. Robots Syst.* IEEE, 2016, pp. 66–73.
- [5] K. Waugh, B. D. Ziebart, and J. A. Bagnell, "Inverse correlated equilibrium for matrix games," *Adv. Neural Inform. Process. Syst.*, 2010.
- [6] V. Kuleshov and O. Schrijvers, "Inverse game theory: Learning utilities in succinct games," in *Int. Conf. Web Internet Econ.* Springer, 2015, pp. 413–427.
- [7] D. Bertsimas, V. Gupta, and I. C. Paschalidis, "Data-driven estimation in equilibrium using inverse optimization," *Math. Prog.*, vol. 153, no. 2, pp. 595–633, 2015.
- [8] L. Peters, D. Fridovich-Keil, V. R. Royo, C. J. Tomlin, and C. Stachniss, "Inferring objectives in continuous dynamic games from noise-corrupted partial state observations," in *Robotics: Sci. and Syst. 2021*, 2021.
- [9] T. L. Molloy, J. I. Charaja, S. Hohmann, and T. Perez, "Inverse optimal control and inverse noncooperative dynamic game theory," 2022.
- [10] Y. Yu, J. Salfity, D. Fridovich-Keil, and U. Topcu, "Inverse matrix games with unique quantal response equilibrium," *IEEE Control Syst. Lett.*, vol. 7, pp. 643–648, 2022.
- [11] J. Li, C.-Y. Chiu, L. Peters, S. Sojoudi, C. Tomlin, and D. Fridovich-Keil, "Cost inference for feedback dynamic games from noisy partial state observations and incomplete trajectories," in *Proc. Int. Conf. Auton. Agents and Multiagent Syst.*, 2023, pp. 1062–1070.
- [12] C. K. Ling, F. Fang, and J. Z. Kolter, "What game are we playing? end-to-end learning in normal and extensive form games," in *Proc. Int. Joint Conf. Artif. Intell.*, 2018, pp. 396–402.
- [13] S. Chen, Y. Yu, D. Fridovich-Keil, and U. Topcu, "Soft-bellman equilibrium in affine markov games: Forward solutions and inverse learning," *arXiv preprint arXiv:2304.00163*, 2023.

- [14] T. Başar, "Affine incentive schemes for stochastic systems with dynamic information," *SIAM J. Control Optim.*, vol. 22, no. 2, pp. 199–210, 1984.
- [15] N. Nisan and A. Ronen, "Algorithmic mechanism design," *Games Econ. Behav.*, vol. 35, no. 1, pp. 359–379, 2015.
- [16] M. Lopes, F. Melo, and L. Montesano, "Active learning for reward estimation in inverse reinforcement learning," in *Proc. Eur. Conf. Mach. Learn. Knowl. Discovery Databases*. Springer Berlin Heidelberg, 2009, pp. 31–46.
- [17] R. Akrou, M. Schoenauer, and M. Sebag, "April: Active preference learning-based reinforcement learning," in *Proc. Eur. Conf. Mach. Learn. Knowl. Discovery Databases*. Springer, 2012, pp. 116–131.
- [18] D. Sadigh, A. Dragan, S. Sastry, and S. Seshia, "Active preference-based learning of reward functions," in *Robotics: Science and Systems XIII*, July 2017.
- [19] P. F. Christiano, J. Leike, T. Brown, M. Martic, S. Legg, and D. Amodei, "Deep reinforcement learning from human preferences," *Adv. Neural Inform. Process. Syst.*, vol. 30, 2017.
- [20] C. Daniel, M. Viering, J. Metz, O. Kroemer, and J. Peters, "Active reward learning," in *Robot.: Sci. Syst.*, vol. 98, 2014.
- [21] E. Bıyık and D. Sadigh, "Batch active preference-based learning of reward functions," in *Proc. Conf. Robot Learn.* PMLR, 2018, pp. 519–528.
- [22] E. Bıyık, D. A. Lazar, D. Sadigh, and R. Pedarsani, "The green choice: Learning and influencing human decisions on shared roads," in *Proc. IEEE Conf. Decision Control*. IEEE, 2019, pp. 347–354.
- [23] M. Palan, G. Shevchuk, N. Charles Landolfi, and D. Sadigh, "Learning reward functions by integrating human demonstrations and preferences," in *Robot.: Sci. Syst.*, 2019.
- [24] S. M. Katz, A.-C. Le Bihan, and M. J. Kochenderfer, "Learning an urban air mobility encounter model from expert preferences," in *Proc. IEEE/AIAA Digit. Avionics Syst. Conf.* IEEE, 2019, pp. 1–8.
- [25] C. Basu, E. Bıyık, Z. He, M. Singhal, and D. Sadigh, "Active learning of reward dynamics from hierarchical queries," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.* IEEE, 2019, pp. 120–127.
- [26] R. Cohn, E. Durfee, and S. Singh, "Comparing action-query strategies in semi-autonomous agents," in *Proc. AAAI Conf. on Artif. Intell.*, vol. 25, no. 1, 2011, pp. 1102–1107.
- [27] C. Daniel, O. Kroemer, M. Viering, J. Metz, and J. Peters, "Active reward learning with a novel acquisition function," *Autonomous Robots*, vol. 39, pp. 389–405, 2015.
- [28] Y. Cui and S. Niekum, "Active reward learning from critiques," in *Proc. IEEE Int. Conf. Robot. Automat.* IEEE, 2018, pp. 6907–6914.
- [29] E. Bıyık, M. Palan, N. C. Landolfi, D. P. Losey, and D. Sadigh, "Asking easy questions: A user-friendly approach to active reward learning," in *Conference on Robot Learning*. PMLR, 2020, pp. 1177–1190.
- [30] V. Myers, E. Bıyık, N. Anari, and D. Sadigh, "Learning multimodal rewards from rankings," in *Proc. Conf. Robot Learn.* PMLR, 2022, pp. 342–352.
- [31] R. Horst and N. V. Thoai, "Dc programming: overview," *J. Optim. Theory Appl.*, vol. 103, pp. 1–43, 1999.
- [32] T. Lipp and S. Boyd, "Variations and extension of the convex–concave procedure," *Optim. Eng.*, vol. 17, pp. 263–287, 2016.
- [33] S. Dempe and A. Zemkoho, "Bilevel optimization," in *Springer Optimization and Its Applications*. Springer, 2020, vol. 161.
- [34] G. Lanckriet and B. K. Sriperumbudur, "On the convergence of the concave-convex procedure," *Adv. Neural Inform. Process. Syst.*, vol. 22, 2009.
- [35] C. M. Bishop and N. M. Nasrabadi, *Pattern recognition and machine learning*. Springer, 2006, vol. 4, no. 4.