

Online inverse reinforcement learning with unknown disturbances

Ryan Self, Moad Abudia and Rushikesh Kamalapurkar

Abstract—This paper addresses the problem of online inverse reinforcement learning for nonlinear systems with modeling uncertainties while in the presence of unknown disturbances. The developed approach observes state and input trajectories for an agent and identifies the unknown reward function online. Sub-optimality introduced in the observed trajectories by the unknown external disturbance is compensated for using a novel model-based inverse reinforcement learning approach. The observer estimates the external disturbances and uses the resulting estimates to learn the dynamic model of the demonstrator. The learned demonstrator model along with the observed suboptimal trajectories are used to implement inverse reinforcement learning. Theoretical guarantees are provided using Lyapunov theory and a simulation example is shown to demonstrate the effectiveness of the proposed technique.

I. INTRODUCTION

Based on the premise that the most succinct representation of the behavior of an entity is its reward structure [1], this paper aims to recover the reward (or cost) function of an agent by observing the agent performing a task and monitoring its state and control trajectories. The reward function estimation is performed in the presence of modeling uncertainties and unknown disturbances. This process of learning an agent’s reward function is known as inverse reinforcement learning (IRL) [1], [2].

IRL methods are proposed in [1] and reward function estimation techniques using IRL for problems formulated as Markov Decision Processes (MDP) are shown in [3]–[5]. Since solutions to the IRL problems are generally not unique, the maximum entropy principle is developed in [6] to help differentiate between the various solutions. In [7], the authors develop a Maximum Causal Entropy IRL technique for infinite time horizon problems where a stationary soft Bellman policy which helps enable the learning of the transition models is utilized. Beyond this, many extensions of IRL include the formulation of feature constructions [8], solving IRL using gradient based methods [9], and game theoretic approaches [10], which suggest the possibility of finding solutions that outperform the expert. IRL is also extended to nonlinear problems using Gaussian Processes, such as [11].

The aforementioned IRL techniques and inverse optimal control methods [12] are extensively utilized to teach autonomous machines to perform specific tasks in an *offline* setting [13]. However, these *offline* approaches do not have the robustness to uncertainties required for online implementation. Inspired by the success of model-based real-time

reinforcement learning methods such as in [14] and [15] and the online IRL/Inverse Optimal Control (IOC) results for linear systems in [16] and [17], this paper presents an IRL technique for nonlinear systems. In this paper, the results of [18] are extended to address the problem of online IRL in the presence of disturbances by developing a recursive IRL technique.

The main contribution of this paper is the development of a novel method for reward function estimation for an agent with unknown dynamics in the presence of disturbances. The developed technique in this paper builds on the previous work in [18] where a batch IRL method is utilized that relies on optimal demonstrations, and as such, does not consider external disturbances affecting the agent being observed. The recursive IRL update results in smoother weight estimates and admits Lyapunov-based performance guarantees. Addressing the complexities resulting for disturbance-induced sub-optimality of the demonstrations is one of the major technical contributions of this paper. The suboptimal observations make model-free IRL methods challenging because they are entirely trajectory driven, and in general, require either optimal or near optimal observations. The novelty in the technique developed in this paper is the use of model-based IRL to compensate for the disturbance-induced sub-optimality. If dynamic models of the agents under observation are unavailable, they need to be learned from data. However, the disturbances make system identification challenging, and the resulting models are typically poor. To overcome this challenge, it is assumed that the observer and demonstrator are co-located and as a result, experience the same disturbance. One can then learn the disturbances using their effects on the observer and use the resulting estimates to learn the dynamic model of the agent under observation. A model-based IRL method can then be deployed to learn the unknown reward function.

The paper is organized as follows: Section II explains the notation used throughout the paper. Section III details the problem formulation and how the additional challenges related to disturbances are addressed. Section IV details the disturbance estimator for this method. Section V shows the developed parameter estimator. Section VI explains the IRL algorithm. Section VII shows a simulation example for the proposed method and Section VIII concludes the paper.

II. NOTATION

The notation \mathbb{R}^n represents the n –dimensional Euclidean space, and the elements of \mathbb{R}^n are interpreted as column vectors, where $(\cdot)^T$ denotes the vector transpose operator. The set of positive integers excluding 0 is denoted by \mathbb{N} . For

The authors are with the School of Mechanical and Aerospace Engineering, Oklahoma State University, Stillwater, OK, USA. {rself, abudia, rushikesh.kamalapurkar}@okstate.edu.

$a \in \mathbb{R}$, $\mathbb{R}_{\geq a}$ denotes the interval $[a, \infty)$, and $\mathbb{R}_{>a}$ denotes the interval (a, ∞) . If $a \in \mathbb{R}^m$ and $b \in \mathbb{R}^n$, then $[a; b]$ denotes the concatenated vector $\begin{bmatrix} a \\ b \end{bmatrix} \in \mathbb{R}^{m+n}$. The notations I_n and 0_n denote the $n \times n$ identity matrix and the zero element of \mathbb{R}^n , respectively. Whenever it is clear from the context, the subscript n is suppressed.

III. PROBLEM FORMULATION

Consider two agents, Agent 1 and Agent 2, where Agent 1 is monitoring the behavior of Agent 2. Agent 1 has the following dynamics

$$\dot{x}_1 = f_1(x_1, u_1) + d_1, \quad (1)$$

where $x_1 : \mathbb{R}_{\geq T_0} \rightarrow \mathbb{R}^n$ is the state, $u_1 : \mathbb{R}_{\geq T_0} \rightarrow \mathbb{R}^m$ is the control, $f_1 : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ are the dynamics, $d_1 : \mathbb{R}_{\geq T_0} \rightarrow \mathbb{R}^n$ is a disturbance acting on Agent 1, and T_0 is the initial time. The dynamics for Agent 2 are

$$\dot{x}_2 = f_2(x_2, u_2) + d_2, \quad (2)$$

where $x_2 : \mathbb{R}_{\geq T_0} \rightarrow \mathbb{R}^n$ is the state, $u_2 : \mathbb{R}_{\geq T_0} \rightarrow \mathbb{R}^m$ is the control, $f_2 : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ are the dynamics, and $d_2 : \mathbb{R}_{\geq T_0} \rightarrow \mathbb{R}^n$ is a disturbance acting on the Agent 2.

Agent 2 is attempting to follow a policy that minimizes the following performance index

$$J(x_0, u(\cdot)) = \int_{T_0}^{\infty} r(x(t; x_0, u(\cdot)), u(t)) dt, \quad (3)$$

where $x(\cdot; x_0, u(\cdot))$ is the trajectory generated by the optimal controller $u(\cdot)$ for the undisturbed dynamics that minimizes the performance index in (3) starting at x_0 and beginning at time T_0 . The main objective of this paper is to estimate the unknown reward function, r , in the presence of disturbances and uncertainties in the dynamics.

The following assumptions are used in the analysis of the paper.

Assumption 1. *The disturbances affecting both agents are identical, i.e. $d_1(t) = d_2(t) = d(t), \forall t$.*

Assumption 2. *The unknown reward function r is quadratic in the control, i.e.,*

$$r(x, u) = Q(x) + u^T R u, \quad (4)$$

where $R \in \mathbb{R}^{m \times m}$ is a positive definite matrix, such that $R = \text{diag}([r_1, \dots, r_m])$, and the function Q can be represented using a neural network as $Q(x) = (W_Q^*)^T \sigma_Q(x) + \epsilon_Q(x)$ is a positive definite function, where $W_Q^* := [q_1, \dots, q_L]^T$ are ideal reward function weights, $\sigma_Q : \mathbb{R}^n \rightarrow \mathbb{R}^L$ are known continuously differentiable features, and $\epsilon_Q : \mathbb{R}^n \rightarrow \mathbb{R}$ is the approximation error.

Assumption 3. *The dynamics for Agent 2 can be expressed as*

$$\dot{x}_2 = f_2^0(x_2, u_2) + \theta_2^T \sigma_2(x_2, u_2) + d, \quad (5)$$

where $f_2^0 : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ denotes the nominal dynamics, $\theta_2^T \sigma_2$ is a parameterized estimate of the uncertain part of

the dynamics, $\theta_2 \in \mathbb{R}^{p \times n}$ is a matrix of unknown constant parameters, and $\sigma_2 : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^p$ are known features.

If Agent 1 and Agent 2 are co-located and of similar size then the disturbances affecting them can be reasonably assumed to be equal. Assumption 2 facilitates the IRL problem formulation in Section VI, and Assumption 3 facilitates the parameter estimation in Section V.

A solution to the two-agent IRL problem described above is proposed in the following.

Due to the unknown disturbance d acting on Agent 2, the observed trajectories corresponding to Agent 2 will no longer be optimal with respect to its unknown performance index. As a result, a purely data-driven implementation of IRL would yield incorrect reward function estimates. Instead, in this paper, the state trajectories for Agent 2 are measured and the reward function is estimated using a model-based approach that compensates for the trajectory deviations. The unknown disturbance, d , is estimated by Agent 1 using its known internal model, and Agent 1 implements a parameter estimator that incorporates the disturbance estimates to calculate the unknown parameters in the dynamics of Agent 2. Finally, both the disturbance and parameter estimates are used by Agent 1 to estimate the unknown reward function that Agent 2 is acting with respect to.

The following sections; disturbance estimation, parameter estimation, and inverse reinforcement learning, are performed in parallel and in real-time.

IV. DISTURBANCE ESTIMATION

While the IRL method discussed in the following can be developed using any disturbance estimator that results in uniform ultimate boundedness of the disturbance estimation error, the following exponential disturbance estimator (inspired by [19]) is used in this paper for ease of exposition. The disturbance estimation is performed only by Agent 1, and for clarity, the subscripts in the dynamics will be omitted in this section.

The unknown disturbance acting on the agents is assumed to be an additive disturbance that is generated from the exogenous linear system

$$\dot{\zeta} = A\zeta, \quad (6)$$

$$d = C\zeta, \quad (7)$$

where $\zeta : \mathbb{R}_{\geq T_0} \rightarrow \mathbb{R}^N$, $A \in \mathbb{R}^{N \times N}$, $C \in \mathbb{R}^{n \times N}$, and $d : \mathbb{R}_{\geq T_0} \rightarrow \mathbb{R}^n$ is the disturbance.

The disturbance estimator is designed as

$$\dot{\hat{\zeta}} = A\hat{\zeta} + K \left(\dot{x} - \left(f(x, u) + \hat{d} \right) \right), \quad (8)$$

and

$$\hat{d} = C\hat{\zeta}, \quad (9)$$

where $K \in \mathbb{R}^{N \times n}$ is a gain matrix.

The following theorem utilizes Lyapunov-based arguments to establish exponential convergence of the disturbance estimator.

Theorem 1. *If $(A - KC)$ is negative definite, then the disturbance estimation error converges exponentially to zero.*

Proof. Define the error for the exogenous linear system as

$$\tilde{\zeta} = \zeta - \hat{\zeta}. \quad (10)$$

Consider the positive definite candidate Lyapunov function

$$V_d(\tilde{\zeta}) = \frac{1}{2} \tilde{\zeta}^T \tilde{\zeta}. \quad (11)$$

Taking the time-derivative of (11), using (6), and (8)

$$\dot{V}_d(\tilde{\zeta}) = \tilde{\zeta}^T \left(A\tilde{\zeta} - A\hat{\zeta} - K \left(\dot{x} - \left(f(x, u) + \hat{d} \right) \right) \right). \quad (12)$$

Using (2) and simplifying, results in

$$\dot{V}_d(\tilde{\zeta}) = \tilde{\zeta}^T \left(A\tilde{\zeta} - K \left(d - \hat{d} \right) \right). \quad (13)$$

Using (7) and (9)

$$\dot{V}_d(\tilde{\zeta}) = \tilde{\zeta}^T (A - KC) \tilde{\zeta}. \quad (14)$$

Using (14), provided $A - KC$ is negative definite, it can be concluded that $\tilde{\zeta}$ converges exponentially to zero. Since $\tilde{d} = C\tilde{\zeta}$, \tilde{d} has the same convergence rate as $\tilde{\zeta}$. \square

V. PARAMETER ESTIMATION

A parameter estimator, motivated by the authors' previous work in [20], is developed in this section. Since parameter estimation is performed only for Agent 2, for clarity, the subscripts for the dynamics will be omitted in this section.

A. Parameter Estimator

Integrating (5) over the interval $[t - T, t]$ for some constant $T \in \mathbb{R}_{>0}$,¹

$$\begin{aligned} x(t) - x(t - T) &= \int_{t-T}^t f^o(x(\gamma), u(\gamma)) d\gamma \\ &+ \theta^T \int_{t-T}^t \sigma(x(\gamma), u(\gamma)) d\gamma + \int_{t-T}^t d(\gamma) d\gamma. \end{aligned} \quad (15)$$

The expression in (15) can be rearranged to form the affine system

$$X(t) = F(t) + \theta^T S(t) + D(t), \quad \forall t \in \mathbb{R}_{\geq T_0} \quad (16)$$

where

$$X(t) := \begin{cases} x(t) - x(t - T), & t \in [T_0 + T, \infty), \\ 0, & t < T_0 + T, \end{cases} \quad (17)$$

$$F(t) := \begin{cases} \int_{t-T}^t f^o(x(\gamma), u(\gamma)) d\gamma, & t \in [T_0 + T, \infty), \\ 0, & t < T_0 + T, \end{cases} \quad (18)$$

¹If the integration interval is selected to be too short, there may not be enough information in the vector X_i to achieve accurate parameter estimation. If the integration interval is selected too long, parameter estimates may not be available during transients where they are needed the most. The development of a reasonable heuristic that guides the selection of the integration interval is a topic for future research.

$$S(t) := \begin{cases} \int_{t-T}^t \sigma(x(\gamma), u(\gamma)) d\gamma, & t \in [T_0 + T, \infty), \\ 0, & t < T_0 + T, \end{cases} \quad (19)$$

and

$$D(t) := \begin{cases} \int_{t-T}^t d(\gamma) d\gamma, & t \in [T_0 + T, \infty), \\ 0, & t < T_0 + T. \end{cases} \quad (20)$$

The affine error system in (16) motivates the adaptive estimation scheme that follows, in which a *concurrent learning* [21] technique is developed that utilizes recorded data stored in a history stack to drive parameter estimation.

A history stack, \mathcal{H}^{PE} , is a set of data points $\left\{ \left(X_i, F_i, S_i, \hat{D}_i \right) \right\}_{i=1}^M$ such that

$$X_i = F_i + \theta^T S_i + \hat{D}_i + \mathcal{E}_i, \quad \forall i \in \{1, \dots, M\}, \quad (21)$$

where $\mathcal{E}_i = D_i - \hat{D}_i$, and

$$\hat{D}(t) := \begin{cases} \int_{t-T}^t \hat{d}(\gamma) d\gamma, & t \in [T_0 + T, \infty), \\ 0, & t < T_0 + T. \end{cases} \quad (22)$$

\mathcal{H}^{PE} is called *full rank* if there exists a constant $\underline{c} \in \mathbb{R}$ such that

$$0 < \underline{c} < \lambda_{\min} \{ \mathcal{S} \}, \quad (23)$$

where the matrix $\mathcal{S} \in \mathbb{R}^{p \times p}$ is defined as $\mathcal{S} := \sum_{i=1}^M S_i S_i^T$. The concurrent learning update law to estimate the unknown parameters is then given by

$$\dot{\hat{\theta}} = \alpha_\theta \Gamma_\theta \sum_{i=1}^M S_i \left(X_i - F_i - \hat{\theta}^T S_i - \hat{D}_i \right)^T, \quad (24)$$

where $\alpha_\theta \in \mathbb{R}_{>0}$ is a constant adaptation gain, and $\Gamma_\theta : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^{p \times p}$ is the least-squares gain updated using the update law

$$\dot{\Gamma}_\theta = \beta_\theta \Gamma_\theta - \alpha_\theta \Gamma_\theta \mathcal{S} \Gamma_\theta, \quad (25)$$

where $\beta_\theta \in \mathbb{R}_{>0}$ is a constant gain. Using arguments similar to [22, Corollary 4.3.2], it can be shown that provided $\lambda_{\min} \{ \Gamma_\theta^{-1}(0) \} > 0$, the least squares gain matrix satisfies

$$\underline{\Gamma}_\theta \mathbf{I}_p \leq \Gamma_\theta(t) \leq \bar{\Gamma}_\theta \mathbf{I}_p, \quad (26)$$

where $\underline{\Gamma}_\theta$ and $\bar{\Gamma}_\theta$ are positive constants, and \mathbf{I}_p denotes an $p \times p$ identity matrix. If a full rank history stack that satisfies (21) is not available a priori, then the data points can be recorded online.

From the Lyapunov analysis in Section V-B, it is observed that the rate of decay for the parameter estimation error is proportional to the minimum singular value of \mathcal{S} . Therefore, to promote faster convergence for the parameter estimates, a minimum singular value maximization algorithm is developed. At each time t , the algorithm takes the current new data point, $(X^*, F^*, S^*, \hat{D}^*)$, and checks if replacing the new data point with any data point currently in the history stack increases the minimum singular value. If the new data point does increase the minimum singular value, that is,

$$\lambda_{\min} \left(\sum_{i \neq j} S_i S_i^T + S_j S_j^T \right) < \frac{\lambda_{\min} \left(\sum_{i \neq j} S_i S_i^T + S^* S^{*T} \right)}{(1 + \psi)}, \quad (27)$$

where λ_{\min} represents the minimum singular value of a matrix and ψ is a positive constant, then the new data point replaces the data point currently in the \mathcal{H}^{PE} that results in the largest increase in the minimum singular value, if not the new point is discarded.

Using Lyapunov arguments, it can be shown (see Section V-B) that the parameter estimation error is directly related to the error \mathcal{E}_i in (21). Using the fact that newer values of \hat{D}_i result in smaller \mathcal{E}_i due to the exponential convergence of the disturbance estimates, a purging algorithm is developed to eliminate inaccurate data from \mathcal{H}^{PE} .

The algorithm maintains two history stacks, a main history stack and a transient history stack, labeled \mathcal{H}^{PE} and \mathcal{G}^{PE} , respectively. As soon as \mathcal{G}^{PE} is full and sufficient time has elapsed since the last purge (see Section V-B), \mathcal{H}^{PE} is emptied and \mathcal{G}^{PE} is copied into \mathcal{H}^{PE} .

B. Analysis

A Lyapunov based analysis, summarized in the following theorem, is performed to show convergence for the parameter estimator developed in Section V-A.

Remark 1. To facilitate the following Lyapunov analysis, the dynamics for the parameter estimation error can be expressed as

$$\dot{\tilde{\theta}} = -\alpha_\theta \Gamma_\theta \mathcal{S} \tilde{\theta} - \alpha_\theta \Gamma_\theta \sum_{i=1}^M S_i \mathcal{E}_i, \quad (28)$$

by using (21) and (24), along with the error being defined as $\tilde{\theta} = \theta - \hat{\theta}$.

The stability result is summarized in the following theorem.

Theorem 2. *Provided the sequences of history stacks, $\mathcal{H}_1^{PE}, \mathcal{H}_2^{PE}, \dots$, are uniformly full rank² and \tilde{d} converges to zero exponentially, then for time intervals $[T_s, T_{s+1}] \forall s \in \mathbb{N}$, as $s \rightarrow \infty$, $\|\tilde{\theta}(T_s)\| \rightarrow 0$.*

Proof. Consider the candidate Lyapunov function

$$V_\theta(\tilde{\theta}, t) = \frac{1}{2} \tilde{\theta}^T \Gamma_\theta^{-1}(t) \tilde{\theta}. \quad (29)$$

Using the bounds in (26), the candidate Lyapunov function satisfies

$$\frac{1}{\bar{\Gamma}_\theta} \|\tilde{\theta}\|^2 \leq V_\theta(\tilde{\theta}, t) \leq \frac{1}{\underline{\Gamma}_\theta} \|\tilde{\theta}\|^2. \quad (30)$$

The time-derivative of (29) results in

$$\dot{V}_\theta(\tilde{\theta}, t) = \tilde{\theta}^T \dot{\Gamma}_\theta^{-1}(t) \tilde{\theta} + \frac{1}{2} \tilde{\theta}^T \dot{\Gamma}_\theta^{-1}(t) \tilde{\theta}. \quad (31)$$

Using (24) and (25), along with the identity $\dot{\Gamma}^{-1} = -\Gamma^{-1} \dot{\Gamma} \Gamma^{-1}$, \dot{V}_θ can be expressed as

$$\dot{V}_\theta(\tilde{\theta}, t) = -\frac{1}{2} \alpha_\theta \tilde{\theta}^T \mathcal{S} \tilde{\theta} - \frac{1}{2} \beta_\theta \tilde{\theta}^T \Gamma_\theta^{-1}(t) \tilde{\theta} - \alpha_\theta \tilde{\theta}^T \sum_{i=1}^M S_i \mathcal{E}_i^T.$$

²The authors definition of uniformly full rank history stacks \mathcal{H}_s^{PE} requires a constant lower bound on \underline{c} in (23) $\forall s \in \mathbb{N}$.

Using the Cauchy-Schwartz inequality, and bounds in (23) and (26), \dot{V}_θ can be bounded by

$$\dot{V}_\theta(\tilde{\theta}, t) \leq -\frac{1}{2} \left(\alpha_\theta \underline{c} + \frac{\beta_\theta}{\bar{\Gamma}_\theta} \right) \|\tilde{\theta}\|^2 + \alpha_\theta \|\tilde{\theta}\| \sum_{i=1}^M \|S_i\| \|\mathcal{E}_i\|. \quad (32)$$

Since the states and controls are bounded, $\|S_i\|$ is bounded for all i . The upper bound is defined as

$$\bar{S} := M \max\{S_i\}, \forall i \in [1, \dots, M]. \quad (33)$$

Using this upper bound and Young's Inequality, \dot{V}_θ becomes

$$\dot{V}_\theta(\tilde{\theta}, t) \leq -A V_\theta(\tilde{\theta}, t) + B, \quad (34)$$

where A and B are defined as

$$A := \frac{\Gamma_\theta}{4} \left(\alpha_\theta \underline{c} + \frac{\beta_\theta}{\bar{\Gamma}_\theta} \right), \quad (35)$$

$$B := \frac{(\alpha_\theta \bar{S} \sum_{i=1}^M \|\mathcal{E}_i\|)^2}{\alpha_\theta \underline{c} + \beta_\theta / \bar{\Gamma}_\theta}. \quad (36)$$

Due to the purging of the data to remove erroneous estimates \hat{d} from \mathcal{H}^{PE} , further analysis is needed to show parameter convergence. Let the purging instances be defined as T_1, T_2, \dots that maintain a minimum dwell time, \mathcal{T} , such that $T_{s+1} - T_s \geq \mathcal{T} > 0, \forall s \in \mathbb{N}$.

Solve equation (34) over any time interval $[T_s, T_{s+1})$, yields

$$\bar{V}_{s+1} \leq \bar{V}_s e^{-A(t-T_s)} + \frac{B_{s+1}}{A}, \quad (37)$$

where $\bar{V}_s \geq \|V_\theta(\tilde{\theta}(T_s), T_s)\|$ and B_{s+1} denotes the value of B over interval $[T_s, T_{s+1})$. Due to the exponentially decreasing error term $\|\mathcal{E}_i\|$, it can be seen that

$$B_s > B_{s+1}, \forall s = 1, 2, \dots, \quad (38)$$

and $\lim_{s \rightarrow \infty} B_s = 0$. Furthermore, the dwell time condition results in the bound

$$\bar{V}_{s+1} \leq \bar{V}_s e^{-A\mathcal{T}} + \frac{B_{s+1}}{A}, \forall s = 0, 1, 2, 3, \dots,$$

If the bounds B_s are selected so that

$$B_{s+1} > 2B_s e^{-A\mathcal{T}}, \forall s = 0, 1, \dots, \quad (39)$$

then

$$\bar{V}_{s+1} \leq \frac{2B_{s+1}}{A}, \forall s = 0, 1, 2, \dots, \quad (40)$$

where $B_0 := \frac{A\bar{V}_0}{2}$. As a result, $\lim_{s \rightarrow \infty} \bar{V}_s = 0$. It can further be concluded that $\|\tilde{\theta}(T_s)\| \rightarrow 0$ as $s \rightarrow \infty$. \square

Remark 2. There is no loss of generality in assuming (39) since the bounds B_s can be artificially inflated to meet (39).

VI. INVERSE REINFORCEMENT LEARNING

The formulation of IRL in the following two sections closely follows the authors' previous work in [18]. In addition, IRL is performed only on Agent 2, and the subscripts for the dynamics are omitted in the next sections.

A. Inverse Bellman Error

Under the premise that Agent 2 implements a feedback controller that would be optimal in a disturbance-free environment, the state and control trajectories, $x(\cdot)$ and $u(\cdot)$, satisfy the Hamilton-Jacobi-Bellman equation

$$H\left(x(t), \nabla_x(V^*(x(t)))^T, u(t)\right) = 0, \forall t \in \mathbb{R}_{\geq 0}, \quad (41)$$

where the unknown optimal value function is $V^* : \mathbb{R}^n \rightarrow \mathbb{R}$ and $H : \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ is the Hamiltonian, defined as $H(x, p, u) := p^T f(x, u) + r(x, u)$. The goal of IRL is to accurately estimate the reward function, r . To aid in the estimation of the reward function, let $\hat{V} : \mathbb{R}^n \times \mathbb{R}^P \rightarrow \mathbb{R}$, $(x, \hat{W}_V) \mapsto \hat{W}_V^T \sigma_V(x) + \epsilon_V(x)$ be a parameterized estimate of the optimal value function V^* , where $\hat{W}_V \in \mathbb{R}^P$ are the estimates of the ideal value function weights W_V^* , $\sigma_V : \mathbb{R}^n \rightarrow \mathbb{R}^P$ are known continuously differentiable features, and $\epsilon_V : \mathbb{R}^n \rightarrow \mathbb{R}$ is the resulting approximation error. Using $\hat{\theta}$, \hat{W}_V , \hat{W}_Q , and \hat{W}_R , which are the estimates of θ , W_V^* , W_Q^* , and $W_R^* := [r_1, \dots, r_m]^T$, respectively, in (41), the inverse Bellman error $\delta' : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^{L+P+m} \times \mathbb{R}^P \rightarrow \mathbb{R}$ is obtained as

$$\delta'\left(x, u, \hat{W}, \hat{\theta}\right) = \hat{W}_V^T \nabla_x \sigma_V(x) \hat{Y}(x, u, \hat{\theta}) + \hat{W}_Q^T \sigma_Q(x) + \hat{W}_R^T \sigma_u(u), \quad (42)$$

where $\sigma_u(u) := [u_1^2, \dots, u_m^2]$ and $\hat{Y}(x, u, \hat{\theta}) = [f^o(x, u) + \hat{g}(x, u, \hat{\theta})]$ where $\hat{g}(x, u, \hat{\theta}) := \hat{\theta}^T \sigma(x, u)$ from (5). Rearranging, (42) becomes

$$\delta'\left(x, u, \hat{W}', \hat{\theta}\right) = \left(\hat{W}'\right)^T \sigma'\left(x, u, \hat{\theta}\right), \quad (43)$$

where $\hat{W}' := [\hat{W}_V; \hat{W}_Q; \hat{W}_R]$ and $\sigma'\left(x, u, \hat{\theta}\right) := [\nabla_x \sigma_V(x) \hat{Y}(x, u, \hat{\theta}); \sigma_Q(x); \sigma_u(u)]$.

B. Inverse Reinforcement Learning Formulation

Using the formulation of the inverse Bellman error in Section VI-A, and control signals, trajectories, and parameter estimates stored in a history stack, denoted as \mathcal{H}^{IRL} , the inverse Bellman error, evaluated at time instances t_1, t_2, \dots, t_N can be formulated into matrix form

$$\Delta' = \hat{\Sigma}' \hat{W}', \quad (44)$$

where $\Delta' := [\delta'_t(t_1); \dots; \delta'_t(t_N)]$, $\delta'_t(t) := \delta'\left(x(t), u(t), \hat{W}', \hat{\theta}(t)\right)$, and $\hat{\Sigma}' := [(\hat{\sigma}'_t)^T(t_1); \dots; (\hat{\sigma}'_t)^T(t_N)]$. The IRL problem can then be solved by finding the solution of the linear system in (44). Since (44) is a homogeneous system of linear equations, it can only be solved up to a scaling factor. Since optimal state and control trajectories are invariant with respect to scaling of the cost function, the scaling ambiguity in (44) is to be expected. Since optimal control behaviours are scale-invariant, there is no loss of generality in resolving the scale ambiguity by assigning a fixed known value to one of the reward function weights.

Taking the first element of \hat{W}_R to be known, the inverse BE in (43) can then be expressed as

$$\delta'\left(x, u, \hat{W}, \hat{\theta}\right) = \hat{W}^T \sigma''\left(x, u, \hat{\theta}\right) + r_1 \sigma_{u1}(u), \quad (45)$$

where $\hat{W} := [\hat{W}_V; \hat{W}_Q; \hat{W}_R^-]$, the vector \hat{W}_R^- denotes \hat{W}_R with the first element removed, $\sigma_{ui}(u)$ denotes the i th element of the vector $\sigma_u(u)$, the vector σ_u^- denotes σ_u with the first element removed, and $\sigma''\left(x, u, \hat{\theta}\right) := [\nabla_x \sigma_V(x) \hat{Y}(x, u, \hat{\theta}); \sigma_Q(x); \sigma_u^-(u)]$.

The closed-form nonlinear optimal controller corresponding to the reward structure in (3) provides the relationship

$$-2Ru(t) = (g'(x(t)))^T (\nabla_x \sigma_V(x(t)))^T W_V^* + (g'(x(t)))^T \nabla_x \epsilon(x(t)), \quad (46)$$

which can be expressed as

$$-2r_1 u_1(t) + \Delta_{u1} = \sigma_{g1} \hat{W}_V \\ \Delta_{u-} = \sigma_g^- \hat{W}_V + 2\text{diag}(u_2, \dots, u_m) \hat{W}_R^-,$$

where $g'(x) := \nabla_u f(x, u)$, σ_{g1} and Δ_{u1} denote the first rows and σ_g^- and Δ_{u-} denote all but the first rows of $\sigma_g(x) := (g'(x))^T (\nabla_x \sigma_V(x))^T$ and $\Delta_u(x) := (g'(x))^T \nabla_x \epsilon(x)$, respectively, and $R^- := \text{diag}([r_2, \dots, r_m])$. For simplification, let $\sigma := \left[\sigma'', \begin{bmatrix} \sigma_g^T \\ \Theta \end{bmatrix} \right]$, where

$$\Theta := \left[0_{m \times 2n}, \begin{bmatrix} 0_{1 \times m-1} \\ 2\text{diag}([u_2, \dots, u_m]) \end{bmatrix} \right]^T.$$

Updating matrix form in (44) by removing the known reward weight results in the linear system

$$-\Sigma_{u1} = \hat{\Sigma} \hat{W} - \Delta', \quad (47)$$

where $\hat{\Sigma} := [\hat{\sigma}_t^T(t_1); \dots; \hat{\sigma}_t^T(t_N)]$, and $\Sigma_{u1} := [\sigma'_{u1}(u(t_1)); \dots; \sigma'_{u1}(u(t_N))]$, where $\hat{\sigma}_t(\tau) := \sigma\left(x(\tau), u(\tau), \hat{\theta}(\tau)\right)$, $\sigma'_{u1}(\tau) := [r_1 \sigma_{u1}(\tau); 2r_1 u_1(\tau); 0_{(m-1) \times 1}]$.

At any time instant t , provided the data stored in the history stack \mathcal{H}^{IRL} satisfies

$$\text{rank}(\hat{\Sigma}) = L + P + m - 1, \quad (48)$$

then the recursive update law

$$\dot{W} = \alpha \Gamma \hat{\Sigma}^T \left(-\hat{\Sigma} W - \Sigma_{u1} \right), \quad (49)$$

is shown to result in UUB estimation of the weights W^* . In (49), $\alpha \in \mathbb{R}_{>0}$ is a constant adaptation gain and $\Gamma : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^{(L+P+m-1) \times (L+P+m-1)}$ is the least-squares gain updated using the update law

$$\dot{\Gamma} = \beta \Gamma - \alpha \Gamma \hat{\Sigma}^T \hat{\Sigma} \Gamma, \quad (50)$$

where $\beta \in \mathbb{R}_{>0}$ is the forgetting factor. Using arguments similar to [22, Corollary 4.3.2], it can be shown that provided $\lambda_{\min} \{\Gamma^{-1}(0)\} > 0$, the least squares gain matrix satisfies

$$\underline{\Gamma}_{L+P+m-1} \leq \Gamma(t) \leq \bar{\Gamma}_{L+P+m-1}, \quad (51)$$

where $\underline{\Gamma}$ and $\bar{\Gamma}$ are positive constants, and I_n denotes an $n \times n$ identity matrix.

C. Analysis

A Lyapunov based analysis is performed to show convergence for the IRL method in Section VI-B.

Definition 1. A sequence of history stacks, $\mathcal{H}_1^{IRL}, \mathcal{H}_2^{IRL}, \dots$, is called uniformly full rank if there exists a constant non-zero lower bound on all of the minimum singular values. More specifically,

$$0 < \underline{\sigma} < \lambda_{\min} \left\{ \hat{\Sigma}_s^T \hat{\Sigma}_s \right\}, \forall s \in \mathbb{N}, \quad (52)$$

where $\underline{\sigma} \in \mathbb{R}_{>0}$.

Remark 3. To facilitate the following Lyapunov analysis, the dynamics for the weight estimation error can be described by

$$\dot{\tilde{W}} = -\alpha \Gamma \hat{\Sigma}^T \left(\hat{\Sigma} \tilde{W} + \Delta_\theta \right), \quad (53)$$

using the fact that $\tilde{W} = W^* - \hat{W}$, along with (49) and the equation $-\Sigma_{u1} = \hat{\Sigma} W^* + \Delta_\theta$, where Δ_θ denotes the errors resulting from poor $\hat{\theta}$ estimates incorporated in $\hat{\Sigma}$.

The stability result is summarized in the following theorem.

Theorem 3. Provided \mathcal{H}^{PE} and \mathcal{H}^{IRL} are uniformly full rank and \hat{d} converges to zero exponentially, then as $t \rightarrow \infty$, $\|\tilde{W}(t)\|$ is uniformly ultimately bounded (UUB).

Proof. Consider the positive definite candidate Lyapunov function

$$V(\tilde{W}, t) = \frac{1}{2} \tilde{W}^T \Gamma^{-1}(t) \tilde{W}. \quad (54)$$

Using the bounds in (51), the candidate Lyapunov function satisfies

$$\underline{v} \|\tilde{W}\|^2 \leq V(\tilde{W}, t) \leq \bar{v} \|\tilde{W}\|^2. \quad (55)$$

where $\underline{v} := 1/2\bar{\Gamma}$ and $\bar{v} := 1/2\underline{\Gamma}$.

The time-derivative of (54) results in

$$\dot{V}(\tilde{W}, t) = \tilde{W}^T \Gamma^{-1}(t) \dot{\tilde{W}} + \frac{1}{2} \tilde{W}^T \dot{\Gamma}^{-1}(t) \tilde{W}. \quad (56)$$

Using (50) and (53), along with the identity $\dot{\Gamma}^{-1} = -\Gamma^{-1} \dot{\Gamma} \Gamma^{-1}$, after simplifying the time-derivative can be expressed as

$$\dot{V}(\tilde{W}, t) = -\frac{1}{2} \alpha \tilde{W}^T \hat{\Sigma}^T \hat{\Sigma} \tilde{W} - \alpha \tilde{W}^T \hat{\Sigma}^T \Delta_\theta - \frac{1}{2} \beta \tilde{W}^T \Gamma^{-1}(t) \tilde{W}.$$

Substituting in $\hat{\Sigma} = \Sigma - \tilde{\Sigma}$

$$\begin{aligned} \dot{V}(\tilde{W}, t) = & -\frac{1}{2} \alpha \tilde{W}^T \hat{\Sigma}^T \hat{\Sigma} \tilde{W} - \alpha \tilde{W}^T \left(\Sigma - \tilde{\Sigma} \right)^T \Delta_\theta \\ & - \frac{1}{2} \beta \tilde{W}^T \Gamma^{-1}(t) \tilde{W}. \end{aligned}$$

Since $\Delta_\theta = \Sigma W^* + \Delta_\epsilon - \tilde{\Sigma} W^*$, substituting in and simplifying yields

$$\begin{aligned} \dot{V}(\tilde{W}, t) = & -\frac{1}{2} \alpha \tilde{W}^T \hat{\Sigma}^T \hat{\Sigma} \tilde{W} - \frac{1}{2} \beta \tilde{W}^T \Gamma^{-1}(t) \tilde{W} \\ & - \alpha \tilde{W}^T \Sigma^T \tilde{\Sigma} W^* + \alpha \tilde{W}^T \hat{\Sigma}^T \tilde{\Sigma} W^* - \alpha \tilde{W}^T \Sigma^T \Delta_\epsilon + \alpha \tilde{W}^T \hat{\Sigma}^T \Delta_\epsilon. \end{aligned}$$

Using the Cauchy-Schwartz inequality, and bounds in (51) and (52), \dot{V} can be bounded by

$$\begin{aligned} \dot{V}(\tilde{W}, t) \leq & -\frac{1}{2} \left(\alpha \underline{\sigma} + \frac{1}{\bar{\Gamma}} \beta \right) \|\tilde{W}\|^2 + \alpha \|\tilde{W}\| \|\Sigma\| \|\tilde{\Sigma}\| \|W^*\| \\ & + \alpha \|\tilde{W}\| \|\tilde{\Sigma}\|^2 \|W^*\| + \alpha \|\tilde{W}\| \|\Sigma\| \|\Delta_\epsilon\| \\ & + \alpha \|\tilde{W}\| \|\tilde{\Sigma}\| \|\Delta_\epsilon\|. \quad (57) \end{aligned}$$

Based on the linearly parameterized reward weights, the norm of the resulting error term Δ_ϵ can be expressed as

$$\bar{\Delta}_\epsilon := \left(N \max_{\substack{x \in x(\cdot) \\ u \in u(\cdot)}} \{ \epsilon_V^2(x) + \epsilon_Q^2(x) + \epsilon_u^2(u) \} \right)^{1/2}.$$

Using this upper bound, \dot{V} becomes

$$\begin{aligned} \dot{V}(\tilde{W}, t) \leq & -\frac{1}{2} \left(\alpha \underline{\sigma} + \frac{1}{\bar{\Gamma}} \beta \right) \|\tilde{W}\|^2 \\ & + \alpha \bar{\Delta}_\epsilon \|\tilde{W}\| \|\Sigma\| + \alpha \bar{\Delta}_\epsilon \|\tilde{W}\| \|\tilde{\Sigma}\| \\ & + \alpha \|\tilde{W}\| \|\Sigma\| \|\tilde{\Sigma}\| \|W^*\| + \alpha \|\tilde{W}\| \|\tilde{\Sigma}\|^2 \|W^*\|. \quad (58) \end{aligned}$$

The term $\|\tilde{\Sigma}\|$ can be expressed in terms of $\tilde{\theta}$ as

$$\|\tilde{\Sigma}\| \leq \|\tilde{\theta}\| \bar{\Sigma}, \quad (59)$$

where

$$\bar{\Sigma} := N \max_{\substack{x \in x(\cdot) \\ u \in u(\cdot)}} \{ \|\nabla_x \sigma_V(x)\| \|\sigma(x, u)\| \}. \quad (60)$$

The term $\|\Sigma\|$, which contains true values of the unknown parameters, is bounded above since it is a function of only true parameters, θ , and bounded states and controls, x and u . Let the upper bound on $\|\Sigma\|$ be denoted as

$$\|\Sigma\| \leq \bar{\Sigma}_\theta. \quad (61)$$

Using (59) and (61), \dot{V} becomes

$$\begin{aligned} \dot{V}(\tilde{W}, t) \leq & -\frac{1}{2} \left(\alpha \underline{\sigma} + \frac{1}{\bar{\Gamma}} \beta \right) \|\tilde{W}\|^2 + \alpha \bar{\Delta}_\epsilon \bar{\Sigma}_\theta \|\tilde{W}\| \\ & + \alpha \bar{\Delta}_\epsilon \bar{\Sigma} \|\tilde{W}\| \|\tilde{\theta}\| + \alpha \bar{\Sigma}_\theta \bar{\Sigma} \|W^*\| \|\tilde{W}\| \|\tilde{\theta}\| \\ & + \alpha \bar{\Sigma}^2 \|W^*\| \|\tilde{W}\| \|\tilde{\theta}\|^2. \quad (62) \end{aligned}$$

Using Young's Inequality \dot{V} then becomes

$$\begin{aligned} \dot{V}(\tilde{W}, t) \leq & -\frac{1}{8} \left(\alpha \underline{\sigma} + \frac{1}{\bar{\Gamma}} \beta \right) \|\tilde{W}\|^2 \\ & + \frac{2 \left(\alpha \bar{\Delta}_\epsilon \bar{\Sigma}_\theta \bar{\theta} + \alpha \bar{\Sigma}_\theta \bar{\Sigma} \|W^*\| \bar{\theta} + \alpha \bar{\Sigma}^2 \|W^*\| \bar{\theta}^2 \right)}{\alpha \underline{\sigma} + \beta / \bar{\Gamma}} \\ & + \frac{2(\alpha \bar{\Delta}_\epsilon \bar{\Sigma} \|\tilde{\theta}\|)^2}{\alpha \underline{\sigma} + \beta / \bar{\Gamma}}, \quad (63) \end{aligned}$$

where $\bar{\theta}$ denotes bounded $\tilde{\theta}$ values stored in the history stack, \mathcal{H}^{IRL} . Using the bound in (55), the differential inequality for \dot{V} can be expressed as

$$\dot{V}(\tilde{W}, t) \leq -AV(\tilde{W}, t) + B + C, \quad (64)$$

where

$$A := \frac{1}{8\bar{v}} \left(\alpha\bar{\sigma} + \frac{1}{\bar{\Gamma}}\beta \right),$$

$$B := \frac{2 \left(\alpha\bar{\Delta}_\epsilon \bar{\Sigma}_\theta \bar{\theta} + \alpha\bar{\Sigma}_\theta \bar{\Sigma} \|W^*\| \bar{\theta} + \alpha\bar{\Sigma}^2 \|W^*\| \bar{\theta}^2 \right)^2}{\alpha\bar{\sigma} + \beta/\bar{\Gamma}}, \quad (65)$$

and

$$C := \frac{2(\alpha\bar{\Delta}_\epsilon \bar{\Sigma} \|\theta\|)^2}{\alpha\bar{\sigma} + \beta/\bar{\Gamma}}. \quad (66)$$

Due to the fact that $\hat{\Sigma}$ and Δ' depend on the quality of the parameter estimates, a purging technique was incorporated in an attempt to remove poor estimates $\hat{\theta}$ from \mathcal{H}^{IRL} . During the transient phase of the parameter estimator, the estimates θ are less accurate and the resulting values of \hat{W} will be poor. Purging facilitates usage of better estimates as they become available.

Due to purging of \mathcal{H}^{IRL} , the estimator is analyzed over discrete time instances. Define the purging instances as T_1, T_2, \dots , and maintain a minimum dwell time, \mathcal{T} , such that $T_{s+1} - T_s \geq \mathcal{T} > 0, \forall s \in \mathbb{N}$.

Solving equation (64) over any time interval $[T_s, T_{s+1})$, yields

$$\bar{V}_{s+1} \leq \bar{V}_s e^{-A(t-T_s)} + \frac{B_s}{A} + \frac{C}{A}, \quad (67)$$

where $\bar{V}_s \geq \left\| V(\tilde{W}(T_s), T_s) \right\|$ and B_{s+1} denotes the value of B over interval $[T_s, T_{s+1})$. A similar argument as the proof of Theorem 2 can be used to conclude that

$$\limsup_{s \rightarrow \infty} \bar{V}_s \leq \frac{32\bar{v}(\alpha\bar{\Delta}_\epsilon \bar{\Sigma} \|\theta\|)^2}{(\alpha\bar{\sigma} + \beta/\bar{\Gamma})^2}, \quad (68)$$

and as a result $\limsup_{s \rightarrow \infty} \left\| \tilde{W}(T_s) \right\| \leq \sqrt{2\frac{\bar{v}}{\bar{v}} \frac{4(\alpha\bar{\Delta}_\epsilon \bar{\Sigma} \|\theta\|)^2}{(\alpha\bar{\sigma} + \beta/\bar{\Gamma})^2}}$. \square

VII. SIMULATION

To demonstrate the performance of the developed method, a nonlinear optimal control problem was constructed using [12] in order to have a known value function for comparison.

Agent 1 has the following nonlinear dynamics

$$\dot{x}_{11} = x_{12}, \quad \dot{x}_{12} = x_{11}x_{12} + 3x_{12}^2 + 5u_1 + d.$$

Agent 2 under observation has the following nonlinear dynamics

$$\begin{aligned} \dot{x}_{21} &= x_{22}, \\ \dot{x}_{22} &= \theta_1 x_{21} \left(\frac{\pi}{2} + \tan^{-1}(5x_{21}) \right) + \frac{\theta_2 x_{21}^2}{1 + 25x_{21}^2} \\ &\quad + \theta_3 x_{22} + 3u_2 + d, \end{aligned} \quad (69)$$

where x_{A_B} denotes state B for Agent A. The parameters θ_1, θ_2 , and θ_3 are unknown constants to be estimated and d is the unknown disturbance. The exact values of these parameters are $\theta_1 = -1, \theta_2 = -\frac{5}{2}$, and $\theta_3 = 4$. The disturbance, d , acting on the agents is generated from the linear system in Section IV, where $A = [0, 1; -1, 0]$ and $C = [0, 0; 1, 0]$, and the chosen gain matrix was $K = [1, 0.5; 0, 5]$.

The performance index that the agent is trying to minimize is

$$J(x_0, u_2(\cdot)) = \int_0^\infty (x_{22}^2 + u_2^2) dt,$$

resulting in the reward function weights to be estimated as $Q = \text{diag}(q_1, q_2) = \text{diag}(0, 1)$ and $R = 1$. The observed state and control trajectories, and the disturbance estimates are used in the estimation of unknown parameters in the dynamics, along with the optimal value function parameters and the reward function weights. The optimal controller is $u_2^* = -3x_{22}$, while the optimal value function is $V^* = x_{21}^2 (v_1 + v_2 \tan^{-1}(5x_{21})) + v_3 x_{22}^2$, resulting in the ideal function parameters $v_1 = \frac{\pi}{2}, v_2 = 1$, and $v_3 = 1$. Figs. 1 and 2

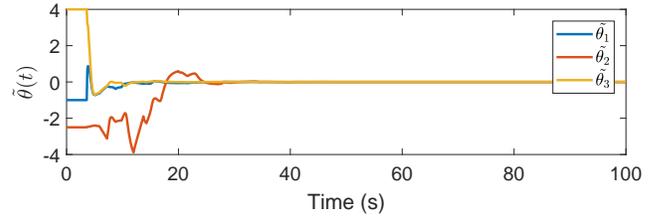


Fig. 1. Estimation error for the unknown parameters in Agent 2's dynamics.

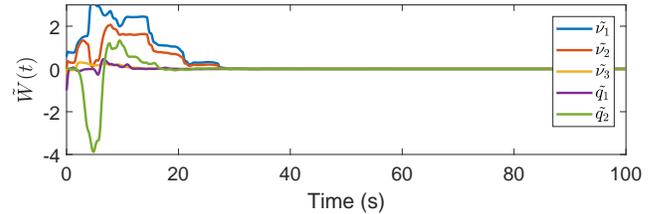


Fig. 2. Estimation error for the unknown parameters in the reward function for Agent 2.

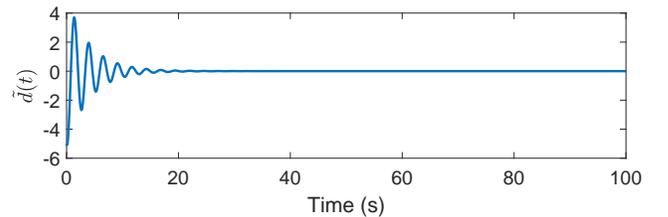


Fig. 3. Estimation error for the unknown disturbance acting on the two agents.

show the performance of the proposed method. Fig. 1 shows convergence of the unknown part of Agent 2's dynamics, and Fig. 2 shows convergence of the unknown reward function. Fig. 3 shows the convergence of the disturbance estimates.

The parameters used for the simulation are: $T = 1.2s$, $N = 100$, $M = 150$, $\beta = \beta_\theta = 0.5$, $\alpha = \alpha_\theta = 1/N$, and a time step of $0.0005s$.

VIII. CONCLUSION

A novel IRL framework is developed in this paper for reward function estimation in the presence of modeling errors and additive disturbances. To compensate for disturbance-induced sub-optimality of observed trajectories, a model-based approach is developed that relies on a disturbance estimator.

Future work will focus on the development of output feedback IRL methods that utilize both state and parameter estimation methods, and extensions of the developed method for disturbances that affect the agents through a control effectiveness matrix. The authors will additionally explore the use of implicit disturbance estimation techniques that would result in bounded disturbance estimation errors.

REFERENCES

- [1] A. Y. Ng and S. Russell, "Algorithms for inverse reinforcement learning," in *Proc. Int. Conf. Mach. Learn.* Morgan Kaufmann, 2000, pp. 663–670.
- [2] S. Russell, "Learning agents for uncertain environments (extended abstract)," in *Proceedings of the eleventh annual conference on Computational learning theory*, 1998.
- [3] P. Abbeel and A. Y. Ng, "Apprenticeship learning via inverse reinforcement learning," in *Proc. Int. Conf. Mach. Learn.*, 2004.
- [4] P. Abbeel and Y. Ng, Andrew, "Exploration and apprenticeship learning in reinforcement learning," in *Proc. Int. Conf. Mach. Learn.*, 2005.
- [5] N. D. Ratliff, J. A. Bagnell, and M. A. Zinkevich, "Maximum margin planning," in *Proc. Int. Conf. Mach. Learn.*, 2006.
- [6] B. D. Ziebart, A. Maas, J. A. Bagnell, and A. K. Dey, "Maximum entropy inverse reinforcement learning," in *Proc. AAAI Conf. Artif. Intel.*, 2008, pp. 1433–1438.
- [7] Z. Zhou, M. Bloem, and N. Bambos, "Infinite time horizon maximum causal entropy inverse reinforcement learning," *IEEE Trans. Autom. Control*, vol. 63, no. 9, pp. 2787–2802, 2018.
- [8] S. Levine, Z. Popovic, and V. Koltun, "Feature construction for inverse reinforcement learning," in *Advances in Neural Information Processing Systems 23*, J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta, Eds. Curran Associates, Inc., 2010, pp. 1342–1350.
- [9] G. Neu and C. Szepesvari, "Apprenticeship learning using inverse reinforcement learning and gradient methods," in *Proc. Anu. Conf. Uncertain. Artif. Intell.* Corvallis, Oregon: AUAI Press, 2007, pp. 295–302.
- [10] U. Syed and R. E. Schapire, "A game-theoretic approach to apprenticeship learning," in *Advances in Neural Information Processing Systems 20*, J. C. Platt, D. Koller, Y. Singer, and S. T. Roweis, Eds. Curran Associates, Inc., 2008, pp. 1449–1456.
- [11] S. Levine, Z. Popovic, and V. Koltun, "Nonlinear inverse reinforcement learning with Gaussian processes," in *Advances in Neural Information Processing Systems 24*, J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2011, pp. 19–27.
- [12] R. E. Kalman, "When is a linear control system optimal?" *J. Basic Eng.*, vol. 86, no. 1, pp. 51–60, 1964.
- [13] K. Mombaur, A. Truong, and J.-P. Laumond, "From human to humanoid locomotion—an inverse optimal control approach," *Auton. Robot.*, vol. 28, no. 3, pp. 369–383, 2010.
- [14] K. G. Vamvoudakis and F. L. Lewis, "Online actor-critic algorithm to solve the continuous-time infinite horizon optimal control problem," *Automatica*, vol. 46, no. 5, pp. 878–888, 2010.
- [15] D. Wang, D. Liu, H. Li, B. Luo, and H. Ma, "An approximate optimal control approach for robust stabilization of a class of discrete-time nonlinear systems with uncertainties," *IEEE Trans. Syst. Man Cybern. Syst.*, vol. 46, no. 5, pp. 713–717, 2016.
- [16] R. Kamalapurkar, "Linear inverse reinforcement learning in continuous time and space," in *Proc. Am. Control Conf.*, Milwaukee, WI, USA, Jun. 2018, pp. 1683–1688.
- [17] T. Molloy, J. Ford, and T. Perez, "Online inverse optimal control on infinite horizons," in *IEEE Conf. Decis. Control.* IEEE, 2018, pp. 1663–1668.
- [18] R. V. Self, M. Harlan, and R. Kamalapurkar, "Online inverse reinforcement learning for nonlinear systems," in *Proc. IEEE Conf. Control Technol. Appl.* Hong Kong, China: IEEE, Aug. 2019, pp. 296–301.
- [19] W.-H. Chen, "Disturbance observer based control for nonlinear systems," *IEEE/ASME Trans. Mechatron.*, vol. 9, no. 4, pp. 706–710, 2004.
- [20] R. Kamalapurkar, "Simultaneous state and parameter estimation for second-order nonlinear systems," in *Proc. IEEE Conf. Decis. Control*, Melbourne, VIC, Australia, Dec. 2017, pp. 2164–2169.
- [21] G. Chowdhary, "Concurrent learning for convergence in adaptive control without persistency of excitation," Ph.D. dissertation, Georgia Institute of Technology, Dec. 2010.
- [22] P. Ioannou and J. Sun, *Robust adaptive control*. Prentice Hall, 1996.