

# Maximum Likelihood Methods for Inverse Learning of Optimal Controllers<sup>\*</sup>

Marcel Menner<sup>\*</sup> and Melanie N. Zeilinger<sup>\*</sup>

<sup>\*</sup> *Institute for Dynamic Systems and Control, ETH Zurich, Zurich, Switzerland (e-mail: {mmenner,mzeilinger}@ethz.ch)*

---

**Abstract:** This paper presents a framework for inverse learning of objective functions for constrained optimal control problems, which is based on the Karush-Kuhn-Tucker (KKT) conditions. We discuss three variants corresponding to different model assumptions and computational complexities. The first method uses a convex relaxation of the KKT conditions and serves as the benchmark. The main contribution of this paper is the proposition of two learning methods that combine the KKT conditions with maximum likelihood estimation. The key benefit of this combination is the systematic treatment of constraints for learning from noisy data with a branch-and-bound algorithm using likelihood arguments. This paper discusses theoretic properties of the learning methods and presents simulation results that highlight the advantages of using the maximum likelihood formulation for learning objective functions.

*Keywords:* Learning for control, data-based control, constrained control.

---

## 1. INTRODUCTION

Objective functions used for control design do not necessarily correspond to the actual performance specifications for a dynamical system, which may comprise complex or sparse targets. Instead, they are often chosen to facilitate gradient-based numerical optimization, which, in turn, makes their design not very intuitive and their calibration can require a tedious manual engineering effort to meet the performance specifications. Inverse learning concepts such as inverse optimal control offer an attractive design paradigm for learning objective functions from data to avoid their manual tuning. In this context, the data can originate, e.g., from a human actor, who demonstrates how to optimally operate the dynamical system being considered. Learning from demonstrations, however, necessarily implies that the data are subject to noise and other sources of sub-optimality, which have to be taken into account.

In this paper, we present and contrast three variants of an inverse optimal control approach that leverages the Karush-Kuhn-Tucker optimality conditions, cf. Karush (1939); Kuhn and Tucker (1951), to learn objective functions of optimal controllers, e.g., for linear quadratic or model predictive control, from noisy data. The first method is based on a convex relaxation of the KKT conditions to allow for noisy data, which is similar to the formulation in Englert et al. (2017); Menner et al. (2019b) and included in this paper as the benchmark. The main contribution of this paper is the proposition of two inverse optimal control methods that combine the KKT conditions with a maximum likelihood estimation algorithm, which offer the key benefit of systematically dealing with state and input constraints in the presence of noisy data. The underlying assumption is that the data are samples

from a distribution (rather than expecting deterministic, optimal data). Maximum likelihood estimation is enabled by an algorithm that uses branch-and-bound-type ideas based on the likelihoods of active constraints. The second contribution is a theoretical and simulative analysis of the properties of the three methods. In theory, we analyze the learning results of the inverse optimal control methods for unconstrained, linear dynamical systems and a quadratic cost function. In simulation, we present learning results for both constrained, linear and nonlinear systems.

*Related work* Inverse optimal control methods typically model data as deterministic and resulting from an optimal control problem [Hewing et al. (2020)], whereas we explicitly consider the data as stochastic. In Kalman (1964), Menner and Zeilinger (2018), and Mombaur et al. (2010), inverse optimal control methods for linear, unconstrained systems are presented. Englert et al. (2017) and Menner et al. (2019b) use a formulation similar to the first method presented in this paper, which is based on the relaxation of the KKT conditions. Chou et al. (2018); Chou et al. (2020) address a related problem by learning constraints. The method in Menner et al. (2019a) considers a non-deterministic model, but does not consider constraints in the learning procedure. The closest to the proposed likelihood estimation methods is Esfahani et al. (2018), where the main difference lies in the proposed formulation using likelihood arguments offering the key advantage of dealing with constraints using a branch-and-bound algorithm.

Inverse reinforcement learning methods typically model data by means of a Markov decision process, cf. Ziebart et al. (2008); Levine and Koltun (2012); Finn et al. (2016). As a result, these methods can deal with noise by construction, but constraints are typically not considered. Compared to inverse reinforcement learning methods, we base our algorithm on the KKT conditions in order to explicitly consider constraints and noisy data.

---

<sup>\*</sup> This work was supported by the Swiss National Science Foundation under grant no. PP00P2\_157601 / 1.

## 2. PROBLEM STATEMENT

We consider discrete-time dynamical systems of the form

$$\mathbf{x}(k+1) = \mathbf{f}(\mathbf{x}(k), \mathbf{u}(k)), \quad (1)$$

where  $\mathbf{x}(k) \in \mathbb{R}^n$  is the state at time  $k$ ,  $\mathbf{u}(k) \in \mathbb{R}^m$  is the input, and  $\mathbf{f}$  is, in general, a nonlinear function.

*Control Model* We consider optimal controllers of the form

$$\begin{aligned} \mathbf{v}_k^* &= \arg \min_{\mathbf{v}_k, \mathbf{z}_k} \theta^T \phi(\mathbf{v}_0, \dots, \mathbf{v}_N, \mathbf{z}_0, \dots, \mathbf{z}_{N+1}) \\ \text{s.t. } \mathbf{z}_{k+1} &= \mathbf{f}(\mathbf{z}_k, \mathbf{v}_k) \quad \forall k = 0, \dots, N \\ \mathbf{g}(\mathbf{z}_k, \mathbf{v}_k) &\leq \mathbf{0} \quad \forall k = 0, \dots, N \\ \mathbf{z}_0 &= \mathbf{z}(0), \end{aligned} \quad (2)$$

where  $\mathbf{v}_k^*$  are optimal inputs at time  $k$ ,  $\mathbf{z}_k$  are the predicted states given inputs  $\mathbf{v}_k$ , and the initial condition is  $\mathbf{z}(0)$ . The minimizers of (2), i.e., the nominal states  $\mathbf{z}_k^*$  and inputs  $\mathbf{v}_k^*$ , express a motion plan and do not necessarily coincide with the measured states  $\mathbf{x}(k)$  and inputs  $\mathbf{u}(k)$ . The objective function is defined by  $\phi$ , which is weighted by the parameters  $\theta$ . The function  $\mathbf{g}$  defines constraints and  $N$  is the prediction horizon. We assume that  $\phi$ ,  $\mathbf{f}$ , and  $\mathbf{g}$  are known and continuously differentiable.

*Assumption on the data* In expectation, the data are assumed to be the solution to (2), i.e., the demonstration, e.g. from a human agent, is modeled as an optimal controller. Due to noise and other sources of sub-optimality, we assume a probability distribution for the data:

- i) We model the initial condition, denoted  $\mathbf{x}(0)$ , as uncertain and assume

$$\mathbf{x}(0) \sim \mathcal{N}(\mathbf{z}(0), \Sigma_0), \quad (3a)$$

i.e.,  $\mathbf{x}(0)$  is Gaussian distributed with mean  $\mathbf{z}(0)$  and covariance  $\Sigma_0$ .

- ii) We model the observed inputs, denoted  $\mathbf{u}(k)$ , as suboptimal and assume

$$\mathbf{u}(k) \sim \mathcal{N}(\mathbf{v}_k^*, \Sigma_k^u) \quad \forall k = 0, \dots, N, \quad (3b)$$

where  $\mathbf{v}_k^*$  are the minimizers of (2).

In the context of learning from data generated by a human agent, eq. (3a) and eq. (3b) model that a human agent may be uncertain about the true initial state  $\mathbf{x}(0)$  and may not execute the intended motion plan optimally.

*Objective* In this paper, we learn the parameters  $\theta$  of the optimal controller in (2) from data represented in the form of state  $\mathbf{x}(k)$  and input  $\mathbf{u}(k)$  measurements satisfying (1) generated, e.g., by a human actor modeled as in (3).

### Notation & Preliminaries

In order to ease exposition, we vectorize sequences

$$\mathbf{U} = \begin{bmatrix} \mathbf{u}(0) \\ \mathbf{u}(1) \\ \vdots \\ \mathbf{u}(N) \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} \mathbf{z}_0 \\ \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_{N+1} \end{bmatrix}, \quad \mathbf{V} = \begin{bmatrix} \mathbf{v}_0 \\ \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_N \end{bmatrix}$$

and use the function  $\mathbf{F}$  relating the vectorized sequences as in (1):  $\mathbf{Z} = \mathbf{F}(\mathbf{V}, \mathbf{z}_0)$ . We use  $\phi(\mathbf{v}_0, \dots, \mathbf{v}_N, \mathbf{z}_0, \dots, \mathbf{z}_{N+1}) = \phi(\mathbf{V}, \mathbf{Z})$ , as well as  $\mathbf{g}(\mathbf{v}_0, \dots, \mathbf{v}_N, \mathbf{z}_0, \dots, \mathbf{z}_{N+1}) = \mathbf{g}(\mathbf{V}, \mathbf{Z})$  equivalently. Further,  $\mathbf{U} \sim \mathcal{N}(\mathbf{V}, \Sigma)$  implies  $\mathbf{u}(k) \sim \mathcal{N}(\mathbf{v}_k^*, \Sigma_k^u)$  for all  $k = 0, \dots, N$ , i.e.,  $\Sigma \in \mathbb{R}^{mN \times mN}$  is block-diagonal with blocks  $\Sigma_k^u$  and we define  $\|\mathbf{x}\|_{\mathbf{X}} = \mathbf{x}^T \mathbf{X} \mathbf{x}$ .

Let  $\text{idx}, \neg\text{idx} \in \{0, 1\}^s$  with  $\text{idx} + \neg\text{idx} = \{1\}^s$ . For a vector  $\boldsymbol{\lambda} \in \mathbb{R}^s$ , we define  $\boldsymbol{\lambda}_{\text{idx}}$  selecting all elements  $\lambda_i$  for which  $\text{idx}_i = 1$  ( $\boldsymbol{\lambda}_{\neg\text{idx}}$  selecting  $\lambda_i$  for which  $\text{idx}_i = 0$ ).  $\boldsymbol{\delta}^i$  is a unit vector with  $\delta_j^i = 1$  if  $i = j$  and  $\delta_j^i = 0$  if  $i \neq j$ .

Consider the optimization problem

$$\begin{aligned} \mathbf{V}^* &= \arg \min_{\mathbf{V}} \theta^T \phi(\mathbf{V}, \mathbf{x}) \\ \text{s.t. } \mathbf{g}(\mathbf{V}, \mathbf{x}) &\leq \mathbf{0}. \end{aligned} \quad (4)$$

The KKT conditions of (4) are given by

$$\text{KKT}_{\theta}(\mathbf{V}^*, \mathbf{x}) = \begin{cases} \nabla_{\mathbf{V}} \mathcal{L}_{\theta}(\mathbf{V}, \mathbf{x})|_{\mathbf{V}=\mathbf{V}^*} = \mathbf{0} \\ \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{V}, \mathbf{F}(\mathbf{V}, \mathbf{x})) = 0 \\ \mathbf{g}(\mathbf{V}, \mathbf{F}(\mathbf{V}, \mathbf{x})) \leq \mathbf{0} \\ \boldsymbol{\lambda} \geq \mathbf{0} \end{cases} \quad (5)$$

with the dual variables  $\boldsymbol{\lambda}$  and the Lagrangian

$$\mathcal{L}_{\theta}(\mathbf{V}, \mathbf{x}) = \theta^T \phi(\mathbf{V}, \mathbf{F}(\mathbf{V}, \mathbf{x})) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{V}, \mathbf{F}(\mathbf{V}, \mathbf{x})).$$

The KKT conditions are necessary for constrained optimization (first-order derivative tests), i.e., any  $\mathbf{V}^*$  locally minimizing (4) satisfies (5). For more details, the reader is referred, e.g., to Boyd and Vandenberghe (2004).

*Proposition 1.* Consider

$$\begin{aligned} f^1 &= \max_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) & f^2 &= \max_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \\ \text{s.t. } g_1(\mathbf{x}) &\leq 0 & \text{s.t. } g_1(\mathbf{x}) &\leq 0, g_2(\mathbf{x}) \leq 0 \end{aligned}$$

and let  $\{\mathbf{x} | g_1(\mathbf{x}) \leq 0\}$  be a non-empty set and  $f(\mathbf{x})$  be bounded. Then,  $f^1 \geq f^2$ .

## 3. INVERSE LEARNING METHODS

This section presents three methods for inverse learning of the objective function, which utilize the KKT conditions in (5) as follows: Suppose  $\mathbf{V}^*$  is the result of (2) for some true parameters  $\theta = \theta_t$  and initial condition  $\mathbf{z}(0)$ . Then,  $\text{KKT}_{\theta}(\mathbf{V}^*, \mathbf{z}(0))$  hold for  $\theta = \theta_t$ . Hence, the KKT conditions can be used to learn  $\theta_t$  given  $\mathbf{V}^*$ . Method 1 is based on a relaxation of the KKT conditions in (5) to allow for noisy data. Methods 2 and 3 are based on maximum likelihood estimation and use the distribution in (3). The three methods vary in computational complexity and model assumptions in the form of approximations.

*Method 1* This method uses a relaxation of the KKT conditions to directly relate the data  $\mathbf{U}, \mathbf{x}(0)$  with  $\theta$ :

$$\begin{aligned} \hat{\theta} &= \arg \min_{\theta, \boldsymbol{\lambda}} \|\nabla_{\mathbf{V}} \mathcal{L}_{\theta}(\mathbf{V}, \mathbf{x}(0))|_{\mathbf{V}=\mathbf{U}}\|_I \\ \text{s.t. } \boldsymbol{\lambda} &\geq \mathbf{0} \\ \boldsymbol{\lambda}_{\neg\text{idx}} &= \mathbf{0}, \end{aligned} \quad (6a)$$

where  $\neg\text{idx}$  indexes inactive constraints,  $g_i(\mathbf{U}, \mathbf{x}(0))$ , with

$$\neg\text{idx}_i = \begin{cases} 1 & \text{if } g_i(\mathbf{U}, \mathbf{x}(0)) < 0 \\ 0 & \text{else.} \end{cases}$$

The main advantage is that (6a) is a convex optimization problem. Compared to (5),  $\nabla_{\mathbf{V}} \mathcal{L}_{\theta}(\mathbf{V}, \mathbf{x}(0))|_{\mathbf{V}=\mathbf{U}} \neq \mathbf{0}$  as well as  $\mathbf{g}(\mathbf{U}, \mathbf{F}(\mathbf{U}, \mathbf{x})) \neq \mathbf{0}$  due to noisy data. Therefore, we minimize  $\nabla_{\mathbf{V}} \mathcal{L}_{\theta}(\mathbf{V}, \mathbf{x}(0))|_{\mathbf{V}=\mathbf{U}}$ , where  $\boldsymbol{\lambda}^T \mathbf{g}(\mathbf{V}, \mathbf{F}(\mathbf{V}, \mathbf{x})) = 0$  with  $\boldsymbol{\lambda} \geq \mathbf{0}$  is relaxed to  $\boldsymbol{\lambda} \geq \mathbf{0}$  and  $\boldsymbol{\lambda}_{\neg\text{idx}} = \mathbf{0}$ .

*Method 2* This method is based on maximum likelihood estimation and uses the expected value of the initial condition in (3a) with  $\mathbf{x}(0) = \mathbf{z}(0)$ . Using the distribution of

the control inputs in (3b), the parameters  $\theta$  are estimated to maximize the probability of observing  $\mathbf{U}$ :

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta} p(\mathbf{V}|\mathbf{U}, \Sigma) \\ \text{s.t. } \mathbf{V} &= \arg \min_{\tilde{\mathbf{V}}} \theta^T \phi(\tilde{\mathbf{V}}, \mathbf{F}(\tilde{\mathbf{V}}, \mathbf{x}(0))) \quad (6b) \\ \text{s.t. } \mathbf{g}(\tilde{\mathbf{V}}, \mathbf{F}(\tilde{\mathbf{V}}, \mathbf{x}(0))) &\leq \mathbf{0}. \end{aligned}$$

*Method 3* This method considers the uncertainty about the initial condition explicitly. The method additionally optimizes over the initial condition with  $\mathbf{x}(0) \sim \mathcal{N}(\mathbf{z}(0), \Sigma_0)$  and the model for the inverse learning problem is given by

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta, \mathbf{z}(0)} p(\mathbf{V}|\mathbf{U}, \Sigma) p(\mathbf{z}(0)|\mathbf{x}(0), \Sigma_0) \\ \text{s.t. } \mathbf{V} &= \arg \min_{\tilde{\mathbf{V}}} \theta^T \phi(\tilde{\mathbf{V}}, \mathbf{F}(\tilde{\mathbf{V}}, \mathbf{z}(0))) \quad (6c) \\ \text{s.t. } \mathbf{g}(\tilde{\mathbf{V}}, \mathbf{F}(\tilde{\mathbf{V}}, \mathbf{z}(0))) &\leq \mathbf{0}. \end{aligned}$$

Both maximum likelihood estimation methods (6b) and (6c) yield bi-level optimization problems that are solved as described in Section 4.

#### 4. ALGORITHM FOR MAXIMUM LIKELIHOOD ESTIMATION

In the following, we outline the algorithm for learning  $\theta$  using Method 2. The algorithm for Method 3 follows analogously. We first replace the likelihood  $p(\mathbf{V}|\mathbf{U}, \Sigma)$  by its logarithmic likelihood  $\log p(\mathbf{V}|\mathbf{U}, \Sigma)$  (log-likelihood) and the lower level optimization problem in (6b) by its KKT conditions in (5). This way, the bi-level optimization problem is replaced by a combinatorial problem due to the complementary slackness condition  $\lambda^T \mathbf{g}(\mathbf{V}, \mathbf{F}(\mathbf{V}, \mathbf{x})) = 0$ :

$$\begin{aligned} p &= \max_{\theta, \mathbf{V}, \lambda, \text{idx}} \log p(\mathbf{V}|\mathbf{U}, \Sigma) \\ \text{s.t. } \text{KKT}_{\theta, \text{idx}}(\mathbf{V}, \mathbf{x}(0)) & \quad (7) \\ \text{idx} &\in \{0, 1\}^s \end{aligned}$$

with  $\text{idx}$  selecting which of the  $s$  constraints are active, i.e.

$$\text{KKT}_{\theta, \text{idx}}(\mathbf{V}, \mathbf{x}) = \begin{cases} \nabla_{\tilde{\mathbf{V}}} \mathcal{L}_{\theta}(\tilde{\mathbf{V}}, \mathbf{x})|_{\tilde{\mathbf{V}}=\mathbf{V}} = \mathbf{0} \\ \mathbf{g}(\mathbf{V}, \mathbf{F}(\mathbf{V}, \mathbf{x})) \leq \mathbf{0} \\ \mathbf{g}(\mathbf{V}, \mathbf{F}(\mathbf{V}, \mathbf{x}))_{\text{idx}} = 0 \\ \lambda \geq \mathbf{0} \\ \lambda_{-\text{idx}} = \mathbf{0}. \end{cases}$$

Solving the combinatorial optimization problem in (7) directly is computationally intensive. However, using likelihood arguments with branch-and-bound-type ideas, (7) becomes practically feasible as outlined in the following.

Algorithm 1 summarizes the procedure to solve (7), which is based on systematically enumerating candidate solutions and is conceptually similar to active-set methods, cf., Murty and Yu (1988). The algorithm aims at reducing the number of times (7) has to be solved for a fixed combination of active constraints, denoted  $\text{idx}^j \in \{0, 1\}^s$ , where we use  $j$  to index the specific combination of active constraints. First (Line 1 in Alg. 1), we solve

$$\begin{aligned} p^0 &= \max_{\theta, \mathbf{V}} \log p(\mathbf{V}|\mathbf{U}, \Sigma) \quad (8) \\ \text{s.t. } \text{KKT}_{\theta, \text{idx}^0=\{0\}^s}(\mathbf{V}, \mathbf{x}(0)), \end{aligned}$$

where  $p^0$  is the log-likelihood of observing  $\mathbf{U}$  and no active constraints ( $\text{idx}^0 = \{0\}^s$ ). Next (Line 2), we compute

an upper bound on the log-likelihood of constraint  $i$ 's activeness given the data  $\mathbf{U}$  as

$$\begin{aligned} \bar{p}^i &= \max_{\mathbf{V}} \log p(\mathbf{V}|\mathbf{U}, \Sigma) \\ \text{s.t. } \mathbf{g}(\mathbf{V}, \mathbf{F}(\mathbf{V}, \mathbf{x}(0))) &\leq \mathbf{0} \quad (9) \\ \mathbf{g}(\mathbf{V}, \mathbf{F}(\mathbf{V}, \mathbf{x}(0)))_{\delta^i} &= \mathbf{0}. \end{aligned}$$

From Proposition 1, constraint  $i$  is not likely to be active if  $\bar{p}^i \leq p^0$  and consequently, constraint  $i$  does not need to be enumerated, which is the first key component of the algorithm's efficiency as the number of possible constraint combinations, denoted  $c$ , can be reduced significantly.

Next (Line 3), we compute upper bounds for the log-likelihood  $p$  in (7) for the fixed combinations of the active constraints  $\text{idx}^j$  (excluding the discarded constraints):

$$\begin{aligned} \bar{p}^j &= \max_{\mathbf{V}} \log p(\mathbf{V}|\mathbf{U}, \Sigma) \\ \text{s.t. } \mathbf{g}(\mathbf{V}, \mathbf{F}(\mathbf{V}, \mathbf{x}(0))) &\leq \mathbf{0} \quad (10) \\ \mathbf{g}(\mathbf{V}, \mathbf{F}(\mathbf{V}, \mathbf{x}(0)))_{\text{idx}^j} &= \mathbf{0} \end{aligned}$$

for all  $j = 1, \dots, c$ . As a result, we obtain  $c$  possible candidate constraint combinations as well as their upper bounds:

$$\mathcal{D} = \{ \{ \text{idx}^1, \bar{p}^1 \}, \{ \text{idx}^2, \bar{p}^2 \}, \dots, \{ \text{idx}^c, \bar{p}^c \} \}. \quad (11)$$

For ease of exposition,  $\mathcal{D}$  is ordered so that  $\bar{p}^j \geq \bar{p}^{j+1}$ . The log-likelihoods  $\bar{p}^j$  are the second key component for the algorithm's efficiency and are used as the stopping criteria, i.e., (7) is solved for  $\text{idx}^j$  starting with  $j = 1$  until  $\bar{p}^j \leq \max\{p^0, \dots, p^{j-1}\}$  (Line 5–9).

---

#### Algorithm 1 Overall algorithm

---

- 1:  $\hat{\theta}, p^0 \leftarrow$  Solve (8) with  $\text{idx} = \{0\}^s$
  - 2: For each  $i = 1, \dots, s$ , compute  $\bar{p}^i$  using (9) and discard constraint  $i$  if  $\bar{p}^i \leq p^0$
  - 3: Compute  $\mathcal{D}$  in (11) using (10)
  - 4:  $j = 1, \hat{p} = p^0$
  - 5: **while**  $\bar{p}^j > \max\{p^0, \dots, p^{j-1}\} = \hat{p}$  **▷** end if less likely
  - 6:  $\theta^j, p^j \leftarrow$  Solve (7) with fixed active constraints  $\text{idx}^j$
  - 7: **if**  $p^j \geq \hat{p}$
  - 8:  $\hat{\theta} \leftarrow \theta^j, \hat{p} \leftarrow p^j$
  - 9:  $j \leftarrow j + 1$
- 

*Remark 1.* We implemented a projected gradient method that uses backtracking line search [Armijo (1966)] to solve both (6a) for Method 1 and (7) with the fixed active constraints  $\text{idx}^j$  for Method 2 and Method 3. Section 6 details the computation times of the three inverse learning methods, which show that the proposed algorithm is computationally feasible.

Fig. 1 illustrates the concept of the upper bounds on the constraint likelihoods (Line 2 in Algorithm 1). In the given example, constraint 4 and constraint 5 do not have to be considered for learning, i.e., do not need to be enumerated. The resulting candidate constraint combinations are

$$\mathcal{D} = \left\{ \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \bar{p}^1 \right\}, \left\{ \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \bar{p}^2 \right\}, \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \bar{p}^3 \right\}, \left\{ \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \bar{p}^4 \right\} \right\}$$

Note that  $\bar{p}^1 = \bar{p}^1$ ,  $\bar{p}^2 = \bar{p}^2$ ,  $\bar{p}^3 = \bar{p}^2$ , and  $\bar{p}^4 = \bar{p}^3$ .

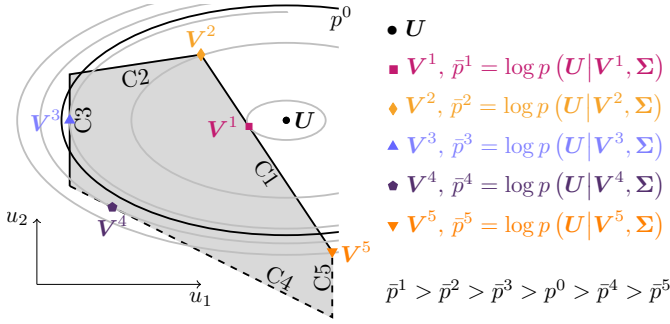


Fig. 1. Illustration of upper bounds on likelihoods  $\bar{p}^i$  as computed in (9) for polytopic constraints  $i = 1, \dots, 5$ . The five upper bounds  $\bar{p}^1 - \bar{p}^5$  (gray ellipses) for the five constraints (C1–C5) as well as  $p^0$  (black ellipse) are displayed as level sets. For  $C_i$ ,  $V^i$  denotes the corresponding, projected input sequence.

## 5. ANALYSIS FOR LINEAR SYSTEMS AND QUADRATIC COST FUNCTION

In this section, we present properties of the learning methods for a common class of dynamical systems and cost functions. Suppose the system in (1) is unconstrained and linear (time-invariant or time-varying), i.e.,

$$f_k(x_k, u_k) = A_k x_k + B_k u_k, \quad (12a)$$

and the cost function is of the form

$$\theta^T \phi(V, Z) = \sum_{k=0}^N z_k^T Q_\theta z_k + v_k^T R_\theta v_k \quad (12b)$$

with  $Q_\theta, R_\theta \succ 0$ . Then, the stationarity condition can be written as

$$\nabla_V \mathcal{L}_\theta(V, z(0)) = M_\theta V + N_\theta z(0),$$

where both  $M_\theta$  and  $N_\theta$  depend on the system dynamics, i.e.,  $A_k, B_k$ , and are linear in their parameters  $\theta$ , i.e.,  $M_{\mu\theta_1 + \theta_2} = \mu M_{\theta_1} + M_{\theta_2}$  and  $N_{\mu\theta_1 + \theta_2} = \mu N_{\theta_1} + N_{\theta_2}$  with the scalar  $\mu$ .

Let  $\theta_t$  be the true parameters and, without loss of generality,  $\|\theta_t\|_2 = 1$  (scale-invariance of the cost function), i.e.,  $\nabla_V \mathcal{L}_{\theta_t}(V, z(0)) = M_{\theta_t} V + N_{\theta_t} z(0) = \mathbf{0}$ . The goal of the learning methods is thus to estimate  $\theta_t$  or any scaled version  $\hat{\theta} = \mu\theta_t$  with  $\mu > 0$  from data. Desirable properties of the learning method are that  $\theta_t$  results in expectation and that all  $\hat{\theta} = \mu\theta_t$  with  $\mu > 0$  are equally likely.

Theorem 1 shows that the expected value of Method 3 is the true parameter vector  $\hat{\theta} = \theta_t$  and that Method 3 is indifferent toward the cost function's scale, i.e., any  $\hat{\theta} = \mu\theta_t$  with  $\mu > 0$  are equally likely (in expectation). Method 2 is equally indifferent toward the parameters' scale but the expected parameters are only  $\hat{\theta} = \theta_t$  if  $x(0) = z(0)$  (proof omitted as it can be similarly derived). Theorem 2 shows that the expected parameters  $\hat{\theta}$  of Method 1 are not necessarily  $\theta_t$  and that Method 1 is not indifferent toward the parameters' scale.

**Theorem 1.** Consider unconstrained, linear systems of the form (12a) and cost functions (12b). Let  $U \sim \mathcal{N}(V, \Sigma)$  and  $x(0) \sim \mathcal{N}(z(0), \Sigma_0)$ . In expectation, Method 3 returns  $\theta_t$  (result 1) and any other parameter realization is necessarily  $\hat{\theta} \propto \theta_t$  (result 2).

**Proof.** Without loss of generality, define  $\theta = \theta_t + \mu\partial\theta$  with  $\mu \in \mathbb{R}$  and  $\partial\theta \in \mathbb{R}^p$  such that  $\|\partial\theta\|_2 = 1$ . The results will be shown by proving the following statements:

*Claim 1:* For any  $\partial\theta$ ,

$$\hat{\mu} = 0 = \arg \max_{\mu, V} \mathbb{E} \left[ -\|U - V\|_{\Sigma^{-1}} - \|x(0) - z(0)\|_{\Sigma_0^{-1}} \right] \quad \text{s.t. } \mathbf{0} = M_{\theta_t + \mu\partial\theta} V + N_{\theta_t + \mu\partial\theta} z(0)$$

*Claim 2:* For  $\partial\theta = \theta_t$ , any  $\mu \in \mathbb{R}$  minimizes

$$\hat{\mu} = \arg \max_{\mu, V} \mathbb{E} \left[ -\|U - V\|_{\Sigma^{-1}} - \|x(0) - z(0)\|_{\Sigma_0^{-1}} \right] \quad \text{s.t. } \mathbf{0} = M_{\theta_t + \mu\partial\theta} V + N_{\theta_t + \mu\partial\theta} z(0)$$

Result 1 follows readily from Claim 1 as  $\theta = \theta_t$ . Result 2 follows from Claim 2 as  $\theta = (1 + \hat{\mu})\theta_t \propto \theta_t$ .

**Proofs of Claim 1 and Claim 2.** Notice first that  $M_\theta \succ \mathbf{0}$  is invertible. The log-likelihood of Method 3 in (6c) is proportional to

$$-\|U - V\|_{\Sigma^{-1}} - \|x(0) - z(0)\|_{\Sigma_0^{-1}}. \quad (13)$$

Using  $V = -M_\theta^{-1} N_\theta z(0)$ , (13) can be written as

$$-\|M_\theta U + N_\theta z(0)\|_{(M_\theta^T \Sigma M_\theta)^{-1}} - \|x(0) - z(0)\|_{\Sigma_0^{-1}}. \quad (14)$$

Then, using linearity ( $M_\theta = M_{\theta_t} + \mu M_{\partial\theta}$  and  $N_\theta = N_{\theta_t} + \mu N_{\partial\theta}$ ),  $U = V + \partial V$ , and  $x(0) = z(0) + \partial z$ , the expected value of (14) yields

$$-\mu^2 \|M_{\partial\theta} V + N_{\partial\theta} z(0)\|_{(M_\theta^T \Sigma M_\theta)^{-1}} - \text{trace}(\Sigma \Sigma^{-1}) - \text{trace}(\Sigma_0 \Sigma_0^{-1}) \quad (15)$$

Therefore, for any  $\partial\theta$ ,  $\mu = 0$  maximizes (15), which proves Claim 1. For  $\partial\theta = \theta_t$ ,  $M_{\partial\theta} V + N_{\partial\theta} z(0) = \mathbf{0}$  and  $\mu \in \mathbb{R}$  minimizes (15), which proves Claim 2.  $\square$

**Theorem 2.** Consider unconstrained linear systems of the form (12a) and cost functions (12b). Let  $U \sim \mathcal{N}(V, \Sigma)$  and  $x(0) \sim \mathcal{N}(z(0), \Sigma_0)$ . Then,  $\theta_t$  does not result in expectation from Method 1 (result 1). Further, Method 1 is not scale-invariant (result 2).

**Proof.** For the considered class of dynamical systems, the cost function in (6a) is

$$\|M_\theta U + N_\theta x(0)\|_I. \quad (16)$$

We define  $\theta = \theta_t + \mu\partial\theta$  with  $\mu \in \mathbb{R}$  and  $\partial\theta \in \mathbb{R}^p$  and  $\|\partial\theta\|_2 = 1$ . Then, using linearity ( $M_\theta = M_{\theta_t} + \mu M_{\partial\theta}$  and  $N_\theta = N_{\theta_t} + \mu N_{\partial\theta}$ ),  $U = V + \partial V$ , and  $x(0) = z(0) + \partial z$ , the expected value of (16) yields

$$\mathbb{E} [\|M_\theta U + N_\theta x(0)\|_I] = \mu^2 (\|M_{\partial\theta} V + N_{\partial\theta} z(0)\|_I + t_{\partial\theta, \partial\theta}) + 2\mu t_{\theta_t, \partial\theta} + t_{\theta_t, \theta_t} \quad (17)$$

with

$$\begin{aligned} t_{\theta_t, \theta_t} &= \text{trace}(M_{\theta_t}^T M_{\theta_t} \Sigma) + \text{trace}(N_{\theta_t}^T N_{\theta_t} \Sigma_0) \\ t_{\theta_t, \partial\theta} &= \text{trace}(M_{\theta_t}^T M_{\partial\theta} \Sigma) + \text{trace}(N_{\theta_t}^T N_{\partial\theta} \Sigma_0) \\ t_{\partial\theta, \partial\theta} &= \text{trace}(M_{\partial\theta}^T M_{\partial\theta} \Sigma) + \text{trace}(N_{\partial\theta}^T N_{\partial\theta} \Sigma_0). \end{aligned}$$

The optimal  $\mu$  minimizing (17) is not necessarily 0 but a function of  $\partial\theta$ , which proves result 1:

$$\mu = -\frac{t_{\theta_t, \partial\theta}}{\|M_{\partial\theta} V + N_{\partial\theta} z(0)\|_I + t_{\partial\theta, \partial\theta}}.$$

For  $\partial\theta = \theta_t$ ,  $\mu = -1$  and  $\theta = \theta_t - \theta_t = \mathbf{0}$ , i.e., the estimator is not scale-invariant, which proves result 2.  $\square$

## 6. SIMULATION RESULTS

In this section, we utilize the three discussed methods for learning the cost function's parameters.

### 6.1 Simulation Setup

*System dynamics and constraints* We consider one linear system and one nonlinear system with dynamics

$$\mathbf{x}(k+1) = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \mathbf{x}(k) + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u(k) \quad (18a)$$

$$\mathbf{x}(k+1) = \begin{bmatrix} 1 & 1 - u(k)^2 \\ 0 & 1 \end{bmatrix} \mathbf{x}(k) + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u(k). \quad (18b)$$

The inputs are constrained as  $|u(k)| \leq 1$ .

*Cost function* The true cost function is chosen as

$$\boldsymbol{\theta}_t^T \boldsymbol{\phi}(\mathbf{V}, \mathbf{Z})$$

with

$$\boldsymbol{\theta}_t = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad \boldsymbol{\phi}(\mathbf{V}, \mathbf{Z}) = \begin{bmatrix} \sum_{k=0}^N z_{1,k}^2 \\ \sum_{k=0}^N z_{2,k}^2 \\ \sum_{k=0}^{N-1} (v_{k+1} - v_k)^2 \\ \sum_{k=0}^N v_k^2 \end{bmatrix}$$

with  $\mathbf{z}_k = [z_{1,k} \ z_{2,k}]^T$  and  $N = 10$ . As the cost function is scale-invariant, we fix one parameter and learn

$$\hat{\boldsymbol{\theta}}^T = [\hat{\theta}_1 \ \hat{\theta}_2 \ \hat{\theta}_3 \ 1].$$

*Data generation* In order to generate the data, we sample the initial conditions  $\mathbf{z}(0)$  from  $\mathbf{z}(0) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Using  $\mathbf{z}(0)$ , we obtain the optimal input sequence  $\mathbf{V}$  using (2). Then, the (sub-optimal) demonstration is generated as  $\mathbf{U} \sim \mathcal{N}(\mathbf{V}, \sigma_u^2 \mathbf{I}_{10})$  and  $\mathbf{x}(0) \sim \mathcal{N}(\mathbf{z}(0), \sigma_0^2 \mathbf{I}_2)$ , where  $\sigma_u$  and  $\sigma_0$  are varied logarithmically as

$$\begin{aligned} \sigma_u, \sigma_0 \in \{ & 0.0001, 0.000215, 0.000464, 0.001, \\ & 0.00215, 0.00464, 0.01, \\ & 0.0215, 0.0464, 0.1, \\ & 0.215, 0.464, 1 \} \end{aligned}$$

*Evaluation criterion* We evaluate the learned parameters by comparing  $\mathbf{V}$  resulting from (2) with  $\hat{\mathbf{V}}$  and  $\hat{\mathbf{V}}$  resulting from (2) with  $\hat{\boldsymbol{\theta}}$  as

$$\text{error} = \frac{\|\hat{\mathbf{V}} - \mathbf{V}\|_2}{\|\mathbf{V}\|_2}. \quad (19)$$

### 6.2 Learning Results

For every tuple  $\{\sigma_u, \sigma_0\}$ , we repeat the data generation process 1000 times and learn  $\hat{\boldsymbol{\theta}}$  using the three inverse learning methods.

Fig. 2 illustrates the median of the error in (19) for the 1000 trials for  $\{\sigma_u, \sigma_0\}$  and the three learning methods. For  $\sigma_u = \sigma_0 = 0$ , error = 0 for all three methods. In the presence of noise  $\sigma_u, \sigma_0 > 0$ , the learning results degrade differently. The error increases for larger standard deviations  $\{\sigma_u, \sigma_0\}$  for all three methods. However, it can be seen that with increased noise levels (sub-optimal data), the error increases more quickly for Method 1, whereas the errors remain smaller for Method 2 and Method 3. Now, consider small  $\sigma_u$ . For increased  $\sigma_0$ , the error for

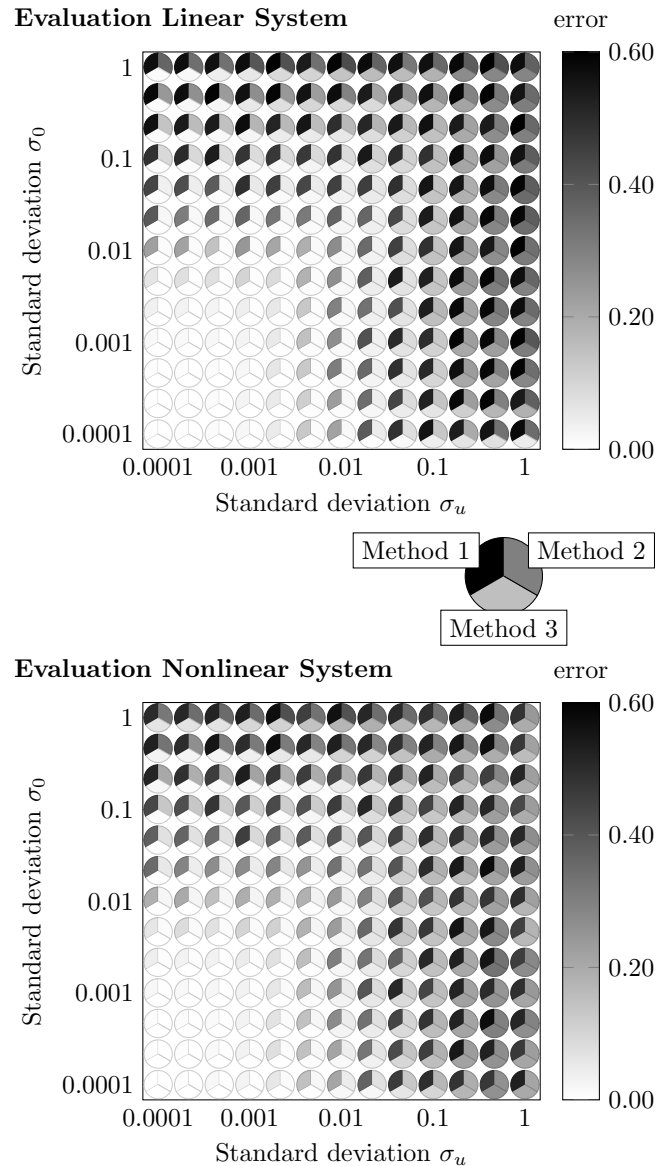
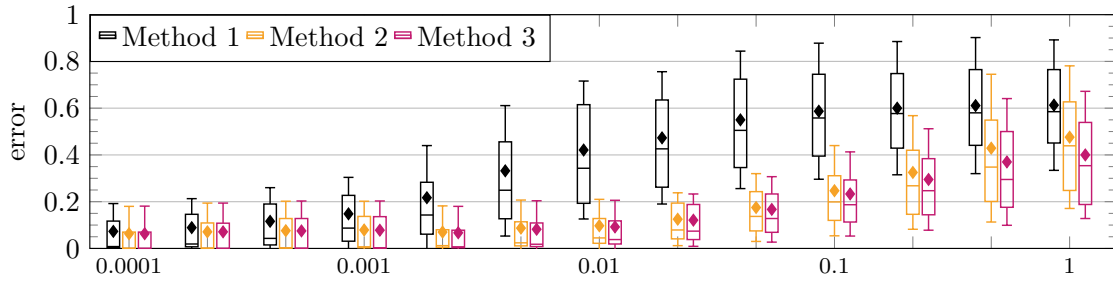


Fig. 2. Median of error  $\in [0, 0.6]$  for different standard deviations  $\sigma_u$  and  $\sigma_0$  (color map from white to black) for the three inverse learning methods.

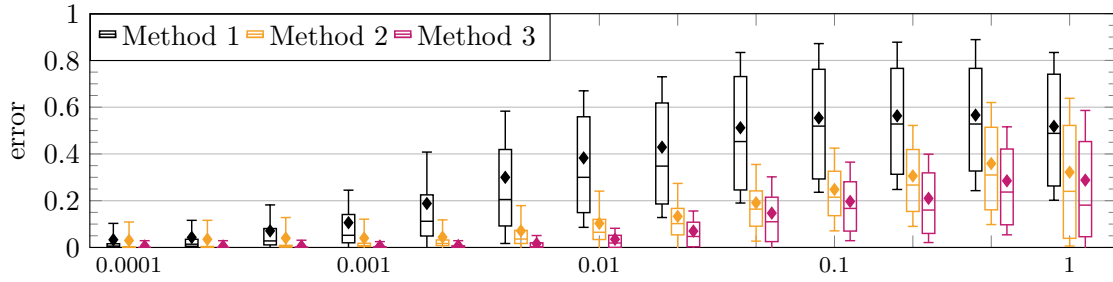
Method 3 is significantly smaller than Method 2, which is expected since Method 3 optimizes over  $\mathbf{z}(0)$ . For small  $\sigma_0$  and increasing  $\sigma_u$ , Method 3 also outperforms Method 2, which suggests that optimizing over the initial condition is advantageous even for small uncertainties in the initial conditions. Note that the standard deviations  $\sigma_u = \sigma_0 > 0.1$  are unrealistically high as  $|u| \leq 1$ .

Fig. 3 shows a more detailed statistical evaluation of the error in (19) for the 1000 trials and  $\sigma_u = \sigma_0$ . First, it can be seen that the estimation with Method 2 and Method 3 have a lower error compared to Method 1. The errors tend to be lower for the nonlinear system compared to the linear system, which can best be seen for  $\sigma_0 = \sigma_u < 0.01$ . The relatively low errors of Method 2 and Method 3 suggest the superiority of the maximum likelihood formulation over the convex relaxation approach of Method 1 measured with respect to the predictive performance, i.e., (19).

### Evaluation Linear System



### Evaluation Nonlinear System



Standard deviation  $\sigma_u, \sigma_0$

Fig. 3. Evaluation of the error in (19) for the three inverse learning methods for different noise levels  $\sigma_u = \sigma_0$ . The diamond symbol and vertical line represent the mean and the median, respectively; the box edges represent the 25th and the 75th percentiles; and the whiskers represent the standard deviation.

### 6.3 Computation Time

Table 1 states the median computation time for all samples of the initial conditions for the three learning methods using MATLAB with the hardware configuration: 3.1 GHz Intel Core i7, 16 GB 1867 MHz DDR3, and Intel Iris Graphics 6100 1536 MB. Method 1 is convex and, therefore, the computation is cheap and requires less than one second. Method 2 is computationally slightly more involved but can still be solved in around one second. Method 3 is more demanding as also the initial condition,  $\mathbf{z}(0)$ , is an optimization variable.

Table 1. Computation time

System	Method	Median over all samples
Linear (18a)	Method 1	$T_{L,M1} = 0.148s$
	Method 2	$T_{L,M2} = 1.19s (8.04T_{L,M1})$
	Method 3	$T_{L,M3} = 3.07s (20.8T_{L,M1})$
Nonlinear (18b)	Method 1	$T_{NL,M1} = 0.142s$
	Method 2	$T_{NL,M2} = 0.584s (4.12T_{NL,M1})$
	Method 3	$T_{NL,M3} = 1.44s (10.2T_{NL,M1})$

Fig. 4 shows a more detailed statistical evaluation of the computation times for  $\sigma_u = \sigma_0$ . The three learning methods have their respective peak computation times at different noise levels, i.e., Method 1's maximum computation times are highest for lower noise levels (peak of median at  $\sigma_u = \sigma_0 = 0.0001$ ); Method 2's maximum times occur for slightly higher noise levels (peak at  $\sigma_u = \sigma_0 = 0.00215$ ); whereas Method 3's peak is for high noise levels (peak at  $\sigma_u = \sigma_0 = 0.0215$ ). For all methods and noise levels, the mean value is higher than the median, which is expected as the median is less susceptible to outliers, i.e., instances with particularly long computation times.

### 7. CONCLUSION

This paper presented three inverse optimal control methods; one method that uses a convex relaxation of the KKT optimality conditions and two methods that combine the KKT conditions with maximum likelihood estimation. It proposed a branch-and-bound-style algorithm for the maximum likelihood formulation, which is based on likelihood arguments to systematically deal with constraints in the presence of noisy data. A simulation study exemplified the three inverse learning methods with both a constrained, linear and a nonlinear system. The results showed that the likelihood estimation methods can be implemented quite efficiently and yield robust learning results, whereas the convex method is computationally efficient but less robust to noise in the training data.

### REFERENCES

Armijo, L. (1966). Minimization of functions having lipschitz continuous first partial derivatives. *Pacific Journal of mathematics*, 16(1), 1–3.

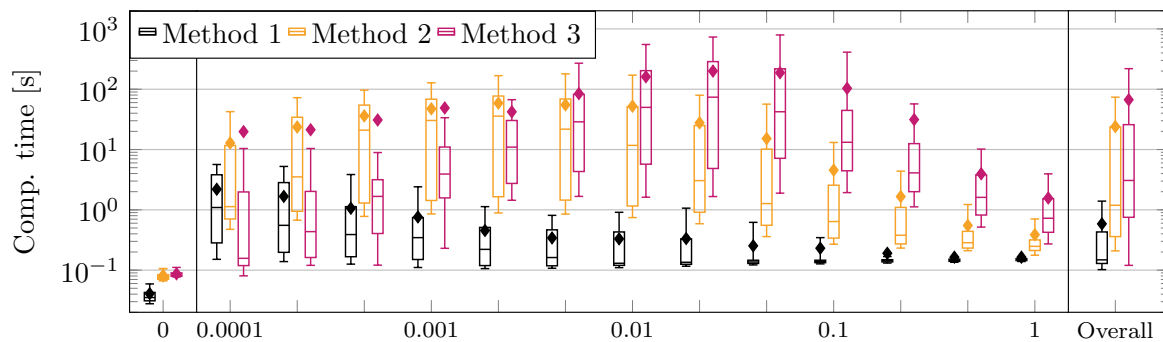
Boyd, S. and Vandenberghe, L. (2004). *Convex optimization*. Cambridge university press.

Chou, G., Ozay, N., and Berenson, D. (2020). Learning constraints from locally-optimal demonstrations under cost function uncertainty. *IEEE Robot. and Automat. Lett.* doi:10.1109/LRA.2020.2974427.

Chou, G., Berenson, D., and Ozay, N. (2018). Learning constraints from demonstrations. *arXiv:1812.07084*.

Englert, P., Vien, N.A., and Toussaint, M. (2017). Inverse KKT: Learning cost functions of manipulation tasks from demonstrations. *Int. J. Robot. Res.*, 36(13–14), 1474–1488.

### Evaluation Linear System



### Evaluation Nonlinear System

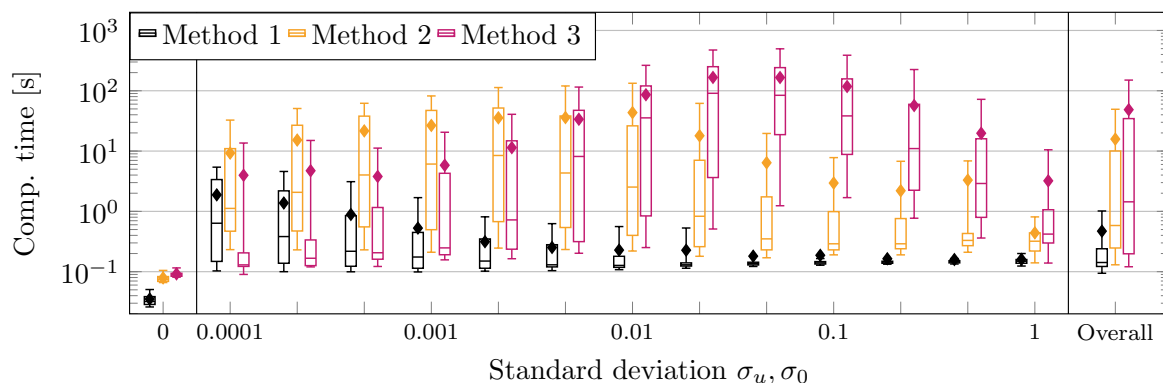


Fig. 4. Statistical evaluation of the computation time for the three inverse learning methods with  $\sigma_u = \sigma_0$ . The diamond symbol and vertical line represent the mean and the median, respectively; the box edges represent the 75th and the 25th percentiles; and the whiskers represent the 90th and the 10th percentiles.

Esfahani, P.M., Shafieezadeh-Abadeh, S., Hanasusanto, G.A., and Kuhn, D. (2018). Data-driven inverse optimization with imperfect information. *Math. Program.*, 167(1), 191–234.

Finn, C., Levine, S., and Abbeel, P. (2016). Guided cost learning: Deep inverse optimal control via policy optimization. In *33rd Int. Conf. Mach. Learn.*, 49–58.

Hewing, L., Wabersich, K.P., Menner, M., and Zeilinger, M.N. (2020). Learning-based model predictive control: Toward safe learning in control. *Annu. Rev. Control, Robot., and Auton. Syst.*, 3, 269–296.

Kalman, R.E. (1964). When is a linear control system optimal? *J. Basic Eng.*, 86(1), 51–60.

Karush, W. (1939). *Minima of functions of several variables with inequalities as side constraints*. Master’s thesis, Dept. of Mathematics, Univ. of Chicago.

Kuhn, H.W. and Tucker, A.W. (1951). Nonlinear programming. In *2nd Berkeley Symposium on Mathematical Statistics and Probability*. University of California Press.

Levine, S. and Koltun, V. (2012). Continuous inverse optimal control with locally optimal examples. In *29th Int. Conf. Machine Learning*, 475–482.

Menner, M., Berntorp, K., Zeilinger, M.N., and Di Cairano, S. (2019a). Inverse learning for human-adaptive motion planning. In *58th IEEE Conf. Decision and Control*, 809–815.

Menner, M., Worsnop, P., and Zeilinger, M.N. (2019b). Constrained inverse optimal control with application to a human manipulation task. *IEEE Trans. Control Syst. Technol.* doi:10.1109/TCST.2019.2955663.

Menner, M. and Zeilinger, M.N. (2018). Convex formulations and algebraic solutions for linear quadratic inverse optimal control problems. In *Eur. Control Conf.*, 2107–2112.

Mombaur, K., Truong, A., and Laumond, J.P. (2010). From human to humanoid locomotion: an inverse optimal control approach. *Auton. Robots*, 28(3), 369–383.

Murty, K.G. and Yu, F.T. (1988). *Linear complementarity, linear and nonlinear programming*, volume 3. Berlin: Heldermann.

Ziebart, B.D., Maas, A., Bagnell, J.A., and Dey, A.K. (2008). Maximum entropy inverse reinforcement learning. In *AAAI Conf. Artif. Intell.*, 1433–1438.