

Inverse KKT – Learning Cost Functions of Manipulation Tasks from Demonstrations

Peter Englert and Marc Toussaint

1 Introduction

Most tasks in real world scenarios require contacts with the environment. For example, the task of opening a door requires contact between the robot gripper and the door handle. In this paper, we address learning from demonstration for the case of manipulation that incorporates contacts. Specifically, we want to extract from demonstrations how to represent and execute manipulations in such a way that the robot can perform such tasks in a robust and general manner.

Cost functions are a powerful representation for robot skills, since they are able to encode task knowledge in a very abstract way. This property allows them to reach high generalization to a wide range of problem configurations. However, designing cost functions by hand can be hard since the right features have to be chosen and combined with each other. Therefore, inverse optimal control, also known as inverse reinforcement learning [18], automates the design of cost functions by extracting the important task spaces and cost parameters from demonstrations. Many successful applications in different areas have demonstrated the capabilities of this idea, including the learning of quadruped locomotion [8], helicopter acrobatics [1] and simulated car driving [10].

There are two parts necessary for applying learning from demonstration with IOC: (1) The inverse optimization method for extracting the cost function from demonstrations; (2) The motion optimization method that creates motions by minimizing such cost functions. Both parts are coupled by the cost function, which is the output of the first and input of the second part, see Fig. 1. Usually IOC algorithms try to find a cost function such that the output of the motion optimization method is similar to

This work was supported by the EU-ICT Project 3rdHand 610878 and the DFG priority program Autonomous Learning 1527.

P. Englert (✉) · M. Toussaint

Machine Learning & Robotics Lab, Universität Stuttgart, Stuttgart, Germany
e-mail: peter.englert@ipvs.uni-stuttgart.de

© Springer International Publishing AG 2018

A. Bicchi and W. Burgard (eds.), *Robotics Research*, Springer Proceedings in Advanced Robotics 3, DOI 10.1007/978-3-319-60916-4_4

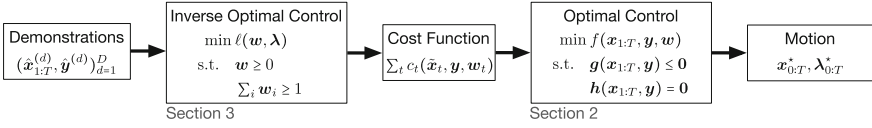


Fig. 1 Concept of skill learning with inverse optimal control, where the cost function plays the central role of encoding the demonstrated behavior. In this paper, we present our formulation of learning a cost function for a constrained trajectory optimization problem

the input demonstrations of the inverse problem. Therefore, the cost function is used as a compact representation that encodes the demonstrated behavior.

Our approach finds a cost function, including the identification of relevant task spaces, such that the demonstrations fulfill the KKT conditions of an underlying constrained optimization problem with this cost function. Thereby we integrate constraints into the IOC method, which allows us to learn from object manipulation demonstrations that naturally involve contact constraints. Motion generation for such cost functions (point 2 above) is a non-linear constrained program, which we solve using an augmented Lagrangian method. However, for typical cost function parameterizations, the IOC problem of inferring the cost function parameters (point 1 above) becomes a quadratic program, which can be solved very efficiently.

The structure of the paper is as follows. We would like to defer the discussion of related work to after we have introduced our method, in Sect. 4. First, in Sect. 2, we introduce some background on constrained trajectory optimization, which represents the counterpart to the IOC approach. Afterwards, we develop our IOC algorithm in Sect. 3 by deriving a cost function based on KKT conditions. In Sect. 5 we evaluate our approach on simulated and real robot experiments.

The main contribution of this paper is the introduction of an IOC method for constrained motions with equality and inequality constraints that is based on the KKT conditions. This method allows to efficiently extract task spaces and parameters from demonstrations.

2 Constrained Trajectory Optimization

We define a trajectory $\mathbf{x}_{0:T}$ as a sequence of $T + 1$ robot configurations $\mathbf{x}_t \in \mathbb{R}^n$. The goal of trajectory optimization is to find a trajectory $\mathbf{x}_{1:T}^*$, given an initial configuration \mathbf{x}_0 , that minimizes a certain objective function

$$f(\mathbf{x}_{1:T}, \mathbf{y}, \mathbf{w}) = \sum_{t=1}^T c_t(\tilde{\mathbf{x}}_t, \mathbf{y}, \mathbf{w}_t). \quad (1)$$

This defines the objective as a sum over cost terms $c_t(\tilde{\mathbf{x}}_t, \mathbf{y}, \mathbf{w}_t)$, where each cost term depends on a k -order tuple of consecutive states $\tilde{\mathbf{x}}_t = (\mathbf{x}_{t-k}, \dots, \mathbf{x}_{t-1}, \mathbf{x}_t)$,

containing the current and k previous robot configurations [22]. This allows us to specify costs on the level of positions, velocities or accelerations (for $k = 2$) in configuration space as well as any task spaces. In addition to the robot configuration state $\tilde{\mathbf{x}}_t$, we use external parameters of the environment \mathbf{y} to contain information that are important for planning the motion (parameters of the environment's configuration, e.g. object positions). These \mathbf{y} usually vary between different problem instances, which is used to generalize the skill to different environment configurations.

We typically assume that the cost terms in Eq. (1) are a weighted sum of squared features,

$$c_t(\tilde{\mathbf{x}}_t, \mathbf{y}, \mathbf{w}_t) = \boldsymbol{\phi}_t(\tilde{\mathbf{x}}_t, \mathbf{y})^\top \text{diag}(\mathbf{w}_t) \boldsymbol{\phi}_t(\tilde{\mathbf{x}}_t, \mathbf{y}), \quad (2)$$

where $\boldsymbol{\phi}_t(\tilde{\mathbf{x}}_t, \mathbf{y})$ are the features and \mathbf{w}_t is the weighting vector at time t . A simple example for a feature is the robot's endeffector position at the end of the motion T relative to the position of a cup. In this example the feature $\boldsymbol{\phi}_T(\tilde{\mathbf{x}}_t, \mathbf{y})$ would compute the difference between the forward kinematics mapping and cup position (given by \mathbf{y}). More complex tasks define body orientations or relative positions between robot and an object. Transition costs are a special type of features, which could be squared torques, squared accelerations or a combination of those, or velocities or accelerations in any task space.

In addition to the task costs we also consider inequality and equality constraints

$$\forall_t \quad \mathbf{g}_t(\tilde{\mathbf{x}}_t, \mathbf{y}) \leq \mathbf{0}, \quad \mathbf{h}_t(\tilde{\mathbf{x}}_t, \mathbf{y}) = \mathbf{0} \quad (3)$$

which are analogous to features $\boldsymbol{\phi}_t(\tilde{\mathbf{x}}_t, \mathbf{y})$ and can refer to arbitrary task spaces. An example for an inequality constraint is the distance to an obstacle, which should not be below a certain threshold. In this example $\mathbf{g}_t(\tilde{\mathbf{x}}_t, \mathbf{y})$ would be the smallest difference between the distance of the robot body to the obstacle and the allowed threshold. The equality constraints are in our approach mostly used to represent persistent contacts with the environment (e.g., \mathbf{h}_t describes the distance between hand and object that should be *exactly* 0). The motivation for using equality constraints for contacts, instead of using cost terms in the objective function as in Eq. (2), is the fact that minimizing costs does not guarantee that they will become 0, which is essential for establishing a contact.

For better readability we transform Eqs. (1) and (3) into vector notation by introducing the vectors \mathbf{w} , $\boldsymbol{\Phi}$, \mathbf{g} and \mathbf{h} that concatenate all elements over time. This allows us to write the objective function of Eq. (1) as

$$f(\mathbf{x}_{1:T}, \mathbf{y}, \mathbf{w}) = \boldsymbol{\Phi}(\mathbf{x}_{1:T}, \mathbf{y})^\top \text{diag}(\mathbf{w}) \boldsymbol{\Phi}(\mathbf{x}_{1:T}, \mathbf{y}) \quad (4)$$

and the overall optimization problem as

$$\begin{aligned}
\mathbf{x}_{1:T}^* &= \underset{\mathbf{x}_{1:T}}{\operatorname{argmin}} f(\mathbf{x}_{1:T}, \mathbf{y}, \mathbf{w}) \\
\text{s.t. } & \mathbf{g}(\mathbf{x}_{1:T}, \mathbf{y}) \leq \mathbf{0} \\
& \mathbf{h}(\mathbf{x}_{1:T}, \mathbf{y}) = \mathbf{0}
\end{aligned} \tag{5}$$

We solve such problems using the augmented Lagrangian method [13]. Therefore, additionally to the solution $\mathbf{x}_{1:T}^*$ we also get the Lagrange parameters $\boldsymbol{\lambda}_{1:T}^*$, which provide information on when the constraints are active during the motion. This knowledge can be used to make the control of interactions with the environment more robust [23]. We use a Gauss–Newton optimization method to solve the unconstrained Lagrangian problem in the inner loop of augmented Lagrangian. For this problem, the gradient is

$$\nabla_{\mathbf{x}_{1:T}} f(\mathbf{x}_{1:T}, \mathbf{y}, \mathbf{w}) = 2\mathbf{J}(\mathbf{x}_{1:T}, \mathbf{y})^\top \operatorname{diag}(\mathbf{w})\boldsymbol{\Phi}(\mathbf{x}_{1:T}, \mathbf{y}) \tag{6}$$

and the Hessian is approximated as in Gauss–Newton as

$$\nabla_{\mathbf{x}_{1:T}}^2 f(\mathbf{x}_{1:T}, \mathbf{y}, \mathbf{w}) \approx 2\mathbf{J}(\mathbf{x}_{1:T}, \mathbf{y})^\top \operatorname{diag}(\mathbf{w})\mathbf{J}(\mathbf{x}_{1:T}, \mathbf{y}), \tag{7}$$

where $\mathbf{J} = \frac{\partial \boldsymbol{\Phi}}{\partial \mathbf{x}}$ is the Jacobian of the features. Using a gradient based trajectory optimization method restricts the class of possible features $\boldsymbol{\Phi}$ to functions that are continuous with respect to \mathbf{x} . However, we will show in the experimental section that this restriction still allows to represent complex behavior like opening a door or sliding a box on a table.

3 Inverse KKT Motion Optimization

We now present an approach to the inverse problem for the constrained trajectory optimization formulation introduced in the previous section. To this end we learn the weight vector \mathbf{w} in Eq. (5) from demonstrations. We assume that D demonstrations of a task are provided with the robot body (e.g., through teleoperation or kinesthetic teaching) and are given in the form $(\hat{\mathbf{x}}_{1:T}^{(d)}, \hat{\mathbf{y}}^{(d)})_{d=1}^D$, where $\hat{\mathbf{x}}_{1:T}^{(d)}$ is the demonstrated trajectory and $\hat{\mathbf{y}}^{(d)}$ is the environment configuration (e.g., object position).

Our IOC objective is derived from the Lagrange function of the problem in Eq. (5)

$$L(\mathbf{x}_{1:T}, \mathbf{y}, \boldsymbol{\lambda}, \mathbf{w}) = f(\mathbf{x}_{1:T}, \mathbf{y}, \mathbf{w}) + \boldsymbol{\lambda}^\top \begin{bmatrix} \mathbf{g}(\mathbf{x}_{1:T}, \mathbf{y}) \\ \mathbf{h}(\mathbf{x}_{1:T}, \mathbf{y}) \end{bmatrix} \tag{8}$$

and the Karush–Kuhn–Tucker (KKT) conditions. The first KKT condition says that for an optimal solution $\mathbf{x}_{1:T}^*$ the condition $\nabla_{\mathbf{x}_{1:T}} L(\mathbf{x}_{1:T}^*, \mathbf{y}, \boldsymbol{\lambda}, \mathbf{w}) = \mathbf{0}$ has to be fulfilled. With Eq. (6) this leads to

$$2\mathbf{J}(\mathbf{x}_{1:T}, \mathbf{y})^\top \text{diag}(\mathbf{w})\Phi(\mathbf{x}_{1:T}, \mathbf{y}) + \boldsymbol{\lambda}^\top \mathbf{J}_c(\mathbf{x}_{1:T}, \mathbf{y}) = \mathbf{0} \quad (9)$$

where the matrix \mathbf{J}_c is the Jacobian of all constraints. We assume that the demonstrations are optimal and should fulfill this conditions. Therefore, the IOC problem can be viewed as searching for a parameter \mathbf{w} such that this condition is fulfilled for all the demonstrations.

We express this idea in terms of the loss function

$$\ell(\mathbf{w}, \boldsymbol{\lambda}) = \sum_{d=1}^D \ell^{(d)}(\mathbf{w}, \boldsymbol{\lambda}^{(d)}) \quad (10)$$

$$\ell^{(d)}(\mathbf{w}, \boldsymbol{\lambda}^{(d)}) = \left(\nabla_{\mathbf{x}_{1:T}} L(\hat{\mathbf{x}}_{1:T}^{(d)}, \hat{\mathbf{y}}^{(d)}, \boldsymbol{\lambda}^{(d)}, \mathbf{w}) \right)^2, \quad (11)$$

where we sum over D demonstrations of the scalar product of the first KKT condition. In Eq. (10), d enumerates the demonstrations and $\boldsymbol{\lambda}^{(d)}$ is the dual to the demonstration $\hat{\mathbf{x}}_{1:T}^{(d)}$ under the problem defined by \mathbf{w} . Note that the dual demonstrations are initially unknown and, of course, depend on the underlying cost function f . More precisely, $\boldsymbol{\lambda}^{(d)} = \boldsymbol{\lambda}^{(d)}(\hat{\mathbf{x}}_{1:T}^{(d)}, \hat{\mathbf{y}}^{(d)}, \mathbf{w})$ is a function of the primal demonstration, the environment configuration of that demonstration, and the underlying parameters \mathbf{w} . And $\ell^{(d)}(\mathbf{w}, \boldsymbol{\lambda}^{(d)}(\mathbf{w})) = \ell^{(d)}(\mathbf{w})$ becomes a function of the parameters only (we think of $\hat{\mathbf{x}}_{1:T}^{(d)}$ and $\hat{\mathbf{y}}^{(d)}$ as given, fixed quantities, as in Eqs. 10 and 11).

Given that we want to minimize $\ell^{(d)}(\mathbf{w})$ we can substitute $\boldsymbol{\lambda}^{(d)}(\mathbf{w})$ for each demonstration by choosing the dual solution that analytically minimizes $\ell^{(d)}(\mathbf{w})$ *subject to* the KKT's complementarity condition

$$\frac{\partial}{\partial \boldsymbol{\lambda}^{(d)}} \ell^{(d)}(\mathbf{w}, \boldsymbol{\lambda}^{(d)}) = \mathbf{0} \quad (12)$$

$$\Rightarrow \boldsymbol{\lambda}^{(d)}(\mathbf{w}) = -(\tilde{\mathbf{J}}_c \tilde{\mathbf{J}}_c^\top)^{-1} \tilde{\mathbf{J}}_c \mathbf{J}^\top \text{diag}(\Phi) \mathbf{w}. \quad (13)$$

Note that here the matrix $\tilde{\mathbf{J}}_c$ is a subset of the full Jacobian of the constraints \mathbf{J}_c that contains only the active constraints during the demonstration, which we can evaluate as \mathbf{g} and \mathbf{h} are independent of \mathbf{w} . This ensures that (13) is the minimizer subject to the complementarity condition. The number of active constraint at each time point has a limit. This limit would be exceeded if more degrees of freedom of the system are constrained than there are available.

By inserting Eq. (13) into Eq. (11) we get

$$\ell^{(d)}(\mathbf{w}) = \underbrace{4\mathbf{w}^\top \text{diag}(\Phi) \mathbf{J} \left(\mathbf{I} - \tilde{\mathbf{J}}_c^\top (\tilde{\mathbf{J}}_c \tilde{\mathbf{J}}_c^\top)^{-1} \tilde{\mathbf{J}}_c \right) \mathbf{J}^\top \text{diag}(\Phi) \mathbf{w}}_{\mathbf{A}^{(d)}} \quad (14)$$

which is the IOC cost per demonstration. Adding up the loss per demonstration and plugging this into Eq. (10) we get a total inverse KKT loss of

$$\ell(\mathbf{w}) = \mathbf{w}^\top \Lambda \mathbf{w} \quad \text{with} \quad \Lambda = 4 \sum_{d=1}^D \Lambda^{(d)}. \quad (15)$$

The resulting optimization problem is

$$\min_{\mathbf{w}} \mathbf{w}^\top \Lambda \mathbf{w} \quad \text{s.t.} \quad \mathbf{w} \geq \mathbf{0} \quad (16)$$

Note that we constrain the parameters \mathbf{w} to be positive. This reflects that we want squared cost features to only positively contribute to the overall cost in Eq. (4).

However, the above formulation may lead to the singular solution $\mathbf{w} = \mathbf{0}$ where zero costs are assigned to all demonstrations, trivially fulfilling the KKT conditions. This calls for a regularization of the problem. In principle there are two ways to regularize the problem to enforce a non-singular solution: First, we can impose positive-definiteness of Eq. (4) at the demonstrations (cf. [10]). Second, as the absolute scaling of Eq. (4) is arbitrary we may additionally add the constraint

$$\begin{aligned} \min_{\mathbf{w}} \quad & \mathbf{w}^\top \Lambda \mathbf{w} \\ \text{s.t.} \quad & \mathbf{w} \geq \mathbf{0}, \quad \sum_i w_i \geq 1 \end{aligned} \quad (17)$$

to our problem formulation (16). We choose the latter option in our experiments.

Equation (17) is a (convex) quadratic program (QP), for which there exist efficient solvers. The gradient $\mathbf{w}^\top \Lambda$ and Hessian Λ are very structured and sparse, which we exploit in our implementations.

In practice we usually use parametrizations on \mathbf{w} . This is useful since in the extreme case, when for each time step a different parameter is used, this leads to a very high dimensional parameter space (e.g., 10 tasks and 300 time steps lead to 3000 parameter). This space can be reduced by using the same weight parameter over all time steps or to activate a task only at some time points. The simplest variant is to use a linear parametrization $\mathbf{w}(\boldsymbol{\theta}) = A\boldsymbol{\theta}$, where $\boldsymbol{\theta}$ are the parameters that the IOC method learns. This parametrization allows a flexible assignment of one parameter to multiple task costs. Further linear parametrizations are radial basis function or B-spline basis functions over time t to more compactly describe smoothly varying cost parameters. For such linear parametrization the problem in Eq. (17) remains a QP that can be solved very efficiently.

Another option we will consider in the evaluations is to use a nonlinear mapping $\mathbf{w}(\boldsymbol{\theta}) = \mathcal{A}(\boldsymbol{\theta})$ to more compactly represent all parameters. For instance, the parameters \mathbf{w} can be of a Gaussian shape (as a function of t), where the mean and variance of the Gaussian is described by $\boldsymbol{\theta}$. Such a parametrization would allow us to learn directly the time point when costs are active. In such a case, the problem is not convex anymore. We address such problems using a general non-linear programming method (again, augmented Lagrangian) and multiple restarts are required with different initializations of the parameter.

Our approach also works in the unconstrained case. In this case the constraint term vanishes in Eq. (9) and the remaining part is the optimality condition of unconstrained optimization, which says that the gradient of the cost function should be equal to zero.

4 Related Work

In the recent years, there has been extensive research on imitation learning and inverse optimal control. In the following section, we will focus on the approaches and methods that are most related to our work of learning cost functions for manipulation tasks. For a broader overview on IOC approaches, we refer the reader to the survey paper of Zhifei and Joo [24] and for an overview on general imitation learning we recommend Argall et al. [3].

4.1 Max-Entropy and Lagrangian-Based IOC Approaches

The work of Levine and Koltun [10] is perhaps the closest to our approach. They use a probabilistic formulation of inverse optimal control that approximates the maximum entropy model of Ziebart et al. [25]. In our framework of trajectory optimization (cf. Sect. 2) this translates to

$$\min_{\mathbf{w}} \nabla_x f^\top (\nabla_x^2 f)^{-1} \nabla_x f - \log |\nabla_x^2 f|. \quad (18)$$

The first term of this equation is similar to our loss in Eq. (10), where the objective is to get small gradients. Additionally, they use the inverse Hessian as a weighting of the gradient. The second term ensures the positive definiteness of the Hessian and also acts as a regularizer on the weights. The learning procedure is performed by maximizing the log-likelihood of the approximated reward function. Instead of enforcing a fully probabilistic formulation, we focus on finite-horizon *constrained* motion optimization formulation with the benefit that it can handle constraints and leads to a fast QP formulation. Further, our formulation also targets at efficiently extracting the relevant task spaces.

Puydupin-Jamin et al. [16] introduced an approach to IOC that also handles linear constraints. It learns the weight parameter \mathbf{w} and Lagrange parameter λ by solving a least-squares optimization problem

$$\min_{\mathbf{w}, \lambda} \left([2J^\top \text{diag}(\Phi) J_c^\top] \begin{bmatrix} \mathbf{w} \\ \lambda \end{bmatrix} + J^{/w} \right)^2 \quad (19)$$

where $^{\wedge w}$ denotes the part in the cost function that is not weighted with \mathbf{w} . The method only addresses equality constraints (no complementarity condition for λ). Our main concern with this formulation is that there are no constraints that ensure that the weight parameter \mathbf{w} do not become 0 or negative. If $\mathbf{J}^{\wedge w}$ is zero, as in our case, the solution is identically zero (\mathbf{w}, λ) . Starting with the KKT condition, they derive a linear residual function that they optimize analytically as the unconstrained least squares. In the experimental section they consider human locomotion with a unicycle model, where they learn one weight parameter of torques and multiple constraints that define the dynamics of the unicycle model and the initial and target position. The idea of using KKT conditions is similar to our approach. However, our formulation allows for inequality constraints and leads to a QP with boundary constraints that ensures that the resulting parameters are feasible. Instead of optimizing for λ , we eliminate λ from the inverse KKT optimization using Eq. (13).

The work of Albrecht et al. [2] learns cost functions for human reaching motions from demonstrations that are a linear combination of different transition types (e.g., jerk, torque). They transformed a bilevel optimization problem, similar to [11], into a constrained optimization problem of the form

$$\min_{\mathbf{x}_{1:T}, \mathbf{w}, \lambda} \left(\phi^{\text{pos}}(\mathbf{x}_T) - \phi^{\text{pos}}(\hat{\mathbf{x}}_T^{(d)}) \right)^2 \quad (20)$$

$$\text{s.t.} \quad \nabla_{\mathbf{x}_{1:T}} L(\mathbf{x}_{1:T}, \mathbf{y}, \lambda, \mathbf{w}) = \mathbf{0} \quad (21)$$

$$\mathbf{h}(\mathbf{x}_{1:T}) = 0 \quad \sum_i w_i = 1 \quad \mathbf{w} \geq 0 \quad (22)$$

The objective is the squared distance between optimal and demonstrated final hand position. They optimize this objective for the trajectory $\mathbf{x}_{1:T}$, the parameter \mathbf{w} and the Lagrange parameter λ with the constraints that the KKT conditions of the trajectory $\mathbf{x}_{1:T}$ are fulfilled. To apply this approach demonstrations are first preprocessed by extracting a characteristic movement with dynamic time warping and a clustering step. Their results show that a combination of different transition costs represent human arm movements best and that they are able to generalize to new hand positions. The advantage of their approach is that they do not only get the parameter weights \mathbf{w} , but also an optimal trajectory $\mathbf{x}_{1:T}^*$ out of the inverse problem in Eqs. (20)–(22). The use of the KKT conditions differs from our approach in two ways. First, they use the KKT conditions in the constrained part of the formulation in Eq. (21), whereas we use them directly as scalar product in the cost function. Second, they use them on the optimization variables $\mathbf{x}_{1:T}$, whereas we use them on the demonstrations $\hat{\mathbf{x}}^{(d)}$ (see Eq. (10)). Instead of minimizing a function directly of the final endeffector position and only learning weights of transition costs, we present a more general solution to imitation learning that can learn transition and task costs in arbitrary feature spaces. Our approach also handles multiple demonstrations directly without preprocessing them to a characteristic movement.

4.2 *Black-Box Inverse Optimal Control*

Black-box optimization approaches are another category of methods for IOC. There, usually an optimization procedure with two layers is used, where in the outer loop black box optimization methods are used to find suitable parameter of the inner motion problem. For this usually no gradients of the outer loop cost function are required.

Mombaur et al. [11] use such a two-layered approach, where they use in the outer loop a derivative free trust region optimization technique and in the inner loop a direct multiple shooting technique. The fitness function of their outer loop is the squared distance between inner loop solution and demonstrations. They apply it on a human locomotion task, where they record demonstration of human locomotion and learn a cost function that they transfer to a humanoid robot. Rückert et al. [17] uses a similar idea to learn movements. They use covariance matrix adaptation [5] in the outer loop to learn policy parameters of a planned movement primitive represented as a cost function. Such methods usually have high computational costs for higher-dimensional spaces since the black box optimizer needs many evaluations. One also needs to find a cost function for the outer loop that leads to reasonable behavior.

Kalakrishnan et al. [7] introduce an inverse formulation of the path integral reinforcement learning method PI^2 [21] to learn objective functions for manipulation tasks. The cost function consists of a control cost and a general state dependent cost term at each time step. They maximize the trajectory likelihood of demonstrations $p(\hat{\mathbf{x}}_{1:T}|\mathbf{w})$ for all demonstrations by creating sampled trajectories around the demonstrations. Further, they L1 regularize \mathbf{w} to only select a subset of the weights. The method is evaluated on grasping tasks.

4.3 *Task Space Extraction*

Jetchev and Toussaint [6] discover task relevant features by training a specific kind of value function, assuming that demonstrations can be modelled as down-hill walks of this function. Similar to our approach, the function is modelled as linear in several potential task spaces, allowing to extract the one most consistent with demonstrations. In Muhlig et al. [12] they automatically select relevant task spaces from demonstrations. Therefore, the demonstrations are mapped on a set of predefined task spaces, which is then searched for the task spaces that best represent the movement. In contrast to these methods, our approach more rigorously extracts task dimensions in the inverse KKT motion optimization framework, including motions that involve contacts.

4.4 Model-Free Imitation Learning

Another approach is the widely used framework of direct imitation learning with movement primitives [14, 15, 19]. They belong to a more direct approach of imitation learning that does not try to estimate the cost function of the demonstration. Instead they represent the demonstrations in a parametrized form that is used to generalize to new situations (e.g., changing duration of motion, adapting the target). Many extensions with different parametrization exist that try to generalize to more complex scenarios [4, 20]. They are very efficient to learn from demonstrations and have been used for manipulation tasks (e.g., pushing a box [9]).

The major difference of such kind of approaches to our method is that they do not need an internal model of the environment, which is sometimes difficult to obtain. However, if such a model is available it can be used to learn a cost function that provide better generalization abilities than movement primitives. This is the case since cost functions are a more abstract representation of task knowledge. Examples of such generalization abilities are demonstrated in Sect. 5 with a box sliding task where we generalized to different box positions and with the door opening task where we generalized to different door angles.

5 Experiments

In the following experimental evaluations, we demonstrate the learning properties and the practical applicability of our approach and compare it to an alternative method in terms of accuracy and learning speed.

For applying an IOC method a set of potential features Φ has to be provided as input. For the following experiments we implemented a simple feature generator to produce a set of potential cost function features in a task independent manner. The used feature types are:

- **Transition features:** Represent the smoothness of the motion (e.g., sum of squared acceleration or torques)
- **Position features:** Represent a body position relative to another body.
- **Orientation features:** Represent orientation of a body relative to another body.

A body is either a part of the robot or belongs to the environment. In the following experiments the time points are either learned with RBF parametrization or they are heuristically extracted from points of interest of the demonstrations (e.g., zero velocity, contact release). We demonstrate in the following experiments that by combining such simple feature types at different time steps into a cost function allows to represent complex behavior.

First, we present on a simple task the ability to reestimate weight functions from optimal demonstrations with different weight parametrizations. Afterwards, we present more complex tasks like sliding a box and opening a door with a real PR2.

5.1 Different Weight Parametrizations in a Benchmark Scenario

The goal of our work is to learn cost functions for finite horizon optimal control problems, including when and how long the costs should be active. In this experiment we test our approach on a simple benchmark scenario. Therefore, we create synthetic demonstrations by optimizing the forward problem with a known ground truth parameter set w^{GT} and test if it is possible to reestimate these parameters from the demonstrations. We create three demonstrations with 50 time steps, where we define that in the time steps 25–30 of these demonstrations the robot endeffector is close to a target position. For this experiments we use a simple robot arm with 7 degree of freedom and the target is a sphere object. We compare the three parametrizations

- **Direct parametrization:** A different parameter is used at each time step (i.e., $w = \theta$) which results in $\theta \in \mathbb{R}^{50}$.
- **Radial basis function:** The basis functions are equally distributed over the time horizon. We use 30 Gaussian basis functions with standard deviation 0.8. This results in $\theta \in \mathbb{R}^{30}$.
- **Nonlinear Gaussian:** A single unnormalized Gaussian weight profile where we have $\theta \in \mathbb{R}^3$ with the weight as linear parameter and the nonlinear parameters are directly the mean and standard deviation. In this case the mean directly corresponds to the time where the activation is highest.

The demonstrations are used as input to our inverse KKT method (see Sect. 3) and the weights are initialized randomly. A comparison of the learned parameters and the ground truth parameter is shown in Fig. 2. The green line represents the ground truth knowledge used for creating the demonstrations. The black dots show the learned parameters of the direct parametrization. The red line shows the learned Gaussian activation and the blue line shows the RBF network. As it can be seen all parametrization detect the right activation region between the time steps 25–30 and approximate the ground truth profile. The Gaussian and RBF parametrization also give some weight to the region outside the actual cost region, which is reasonable since in the demonstrations the robot is still close to the target position. After learning with these parametrizations, we conclude that the linear RBF network are most suited to learn time profiles of cost functions. The main reason for this is the linearity of the para-

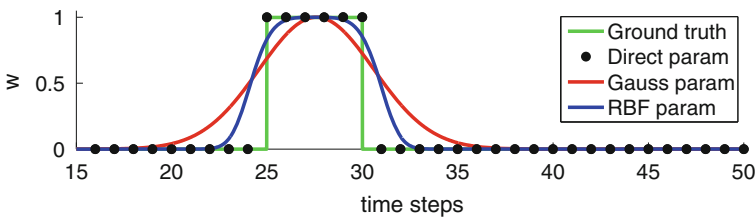


Fig. 2 Learned time profiles of different weight parameterizations. For more details see Sect. 5.1

metrization that makes the inverse KKT problem convex and the versatility of the RBF network to take on more complex forms. Directly learning the time with the nonlinear Gaussian-shaped parametrization was more difficult and required multiple restarts with different initialization. This demonstrates that the framework of constrained trajectory optimization and its counterpart inverse KKT works quite well for reestimating cost functions of optimal demonstrations.

5.2 Sliding a Box on a Table

In this experiment we use our approach to learn a cost function for sliding a box on a table. This task is depicted in Fig. 3. The goal is to move the blue box on the table to the green marked target position and orientation. The robot consist of a fixed base and a hand with 2 fingers. In total the robot has 10 degrees of freedom. Additionally to these degree of freedom we model the box as part of the configuration state, which adds 3 more degrees of freedom (2 translational + 1 rotational). The final box position and orientation is provided as input to our approach and part of the external parameters \mathbf{y} . We used three synthetic demonstrations of the task and created a set of features with the approach described above that led to $\theta \in \mathbb{R}^{537}$ parameters. The relevant features extracted from our algorithm are

- **transition**: Squared acceleration at each time step in joint space
- **posBox**: Relative position between the box and the target.
- **vecBox**: Relative orientation between the box and the target.
- **posFinger1/2**: Relative position between the robots fingertips and the box.
- **posHand**: Relative position between robot hand and box.
- **vecHand**: Relative orientation between robot hand and box.

The contacts between the fingers and the box during the sliding are modeled with equality constraints. They ensure that during the sliding the contact is maintained. For achieving realistic motions, we use an inequality constraint that restrict the movement direction during contact into the direction in which the contact is applied. This ensures that no unrealistic motions like sliding backwards or sideways are created. For clarity we would like to note that we are not doing a physical simulation of the sliding

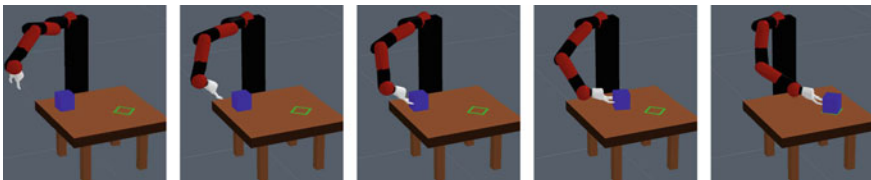


Fig. 3 These images show the box sliding motion of Sect. 5.2 where the goal of the task is to slide the *blue* box on the table to the *green* target region

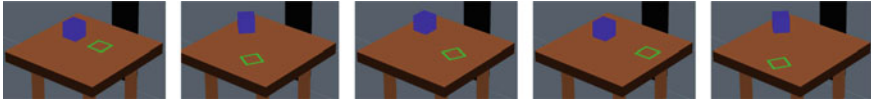


Fig. 4 Each image shows a different instance of the box sliding task. We were able to generalize to different initial box states (*blue* box) and to different final box targets (*green* area)

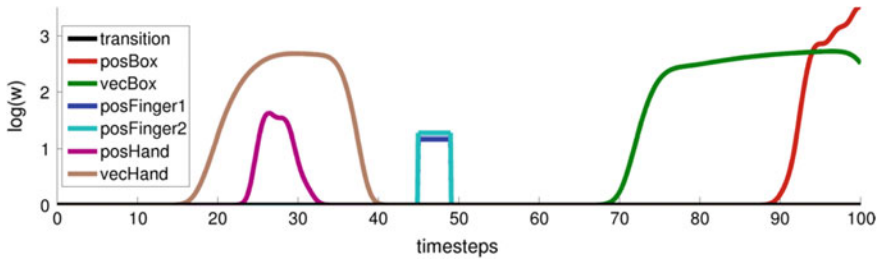


Fig. 5 The resulting parameters w of the extracted relevant features plotted over time. task is depicted in this slideshow

Black box IOC:

repeat

Resample parameters $\{w^{(n)}\}_{n=1}^N$ with CMA

for all $w^{(n)}$ do

Optimize cost function with parameter $w^{(n)}$

Compute fitness $f^{(n)} = \sum_d (x^{(n)} - \hat{x}^{(d)})^2$

end for

Update CMA distribution with fitness values

until

Method	$(x^{(n)} - \hat{x})^2$	comp. time
inverse KKT	0.00021	49.29 sec
black box IOC	0.00542	7116.74 sec

Fig. 6 On the *left* side is the black box IOC algorithm we used for comparison in Sect. 5.2. On the *right* side are the results of the evaluation that show that our method is superior in terms of squared error between the trajectories and computation time

behavior in these experiments. Our goal was more to learn a policy that executes a geometric realistic trajectory from an initial to a final box position. Figure 3 shows one of the resulting motion after learning. We were able to generalize to a wide range of different start and goal position of the box (see Figs. 4 and 5). Videos of the resulting motions can be found in the supplementary material.

We compare our method to a black-box optimization approach similar to [11, 17]. We implemented this approach with the black-box method Covariance Matrix Adaptation (CMA) [5] in the outer loop and our constrained trajectory optimization method (see Sect. 2) in the inner loop. The resulting algorithm is described in Fig. 6. As fitness function for CMA we used the squared distance between the current solution $x^{(n)}$ and the demonstrations $\hat{x}^{(d)}$. We compare this method with our inverse KKT approach by computing the error between the solution and demonstrations and the computational time, which are shown in the table in Fig. 6. The black-box method took around 4900 iterations of the outer loop of the above algorithm until it converged



Fig. 7 The resulting motion after learning the door opening task is depicted in this slideshow. See Sect. 5.3 for more details

to a solution. This comparison shows that using structure and optimality conditions of the solution can enormously improve the learning speed. Further difficulties with black box methods is that they cannot naturally deal with constraints (in our case $w > 0$) and that the initialization is non-trivial.

5.3 Opening a Door with a PR2

In this experiment we apply the introduced inverse KKT approach from Sect. 3 on a skill with the goal to open a door with a real PR2 robot. The problem setup is visualized in Fig. 7. We use a model of the door for our planning approach and track the door angle with AR marker. We use the left arm of the robot that consists of 7 rotational joints and also include the door angle as configuration state into \mathbf{x} . This allows us to define cost functions directly on the door angle. The gripper is fixed during the whole motion. For our IOC algorithm we recorded 2 demonstrations of opening the door from different initial positions with kinesthetic teaching. The motions also include the unlocking of the door by turning the handle first. During the demonstrations we also recorded the door position with the attached markers. We created a feature set similar to the box sliding motion from the previous experiment. Our inverse KKT algorithm extracted the features:

- Relative position & orientation between gripper and handle before and after unlocking the handle.
- Endeffector orientation during the whole opening motion.
- Position of the final door state.

We use equality constraints, similar to the box sliding experiment to keep the contact between endeffector and door. Furthermore, we use inequality constraints to avoid contacts with the rest of the robot body. A resulting motion of optimizing the constrained trajectory optimization problem with the learned parameter \mathbf{w}^* is visualized in Fig. 7. We are able to robustly generate motions with these parameters that generalize to different initial positions and different target door angles (see Fig. 8). Videos of all these motions can be found in the supplementary material.

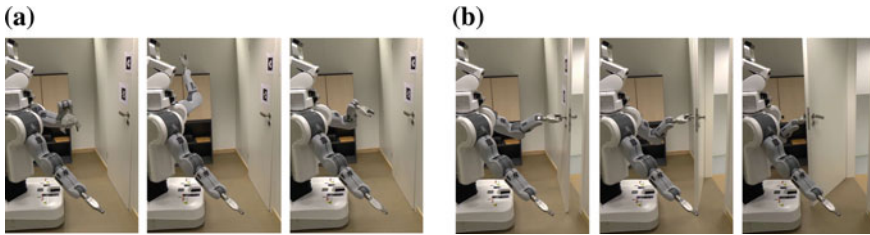


Fig. 8 These images show the generalization abilities of our approach. The pictures in (a) show different initial positions of the robot and the pictures in (b) show different final door angle positions. After learning the weight parameter \mathbf{w}^* with inverse KKT it was possible to generalize to all these instances of the door opening task

6 Conclusion

In this paper we introduced inverse KKT motion optimization, an inverse optimal control method for learning cost functions for constrained motion optimization problems. Our formulation is focused on finite horizon optimal control problems for tasks that include contact with the environment. The resulting method is based on the KKT conditions that the demonstrations should fulfill. For a typical linear parameterization of cost functions this leads to a convex problem; in the general case it is implemented as a 2nd order optimization problem, which leads to a fast convergence rate. We demonstrated the method in a real robot experiment of opening a door that involved contact with the environment. In our future research we plan to further automate and simplify the skill acquisition process. Thereby, one goal is to extend the proposed method to be able to handle demonstrations that are not recorded on the robot body. Another goal is to further improve the skill with reinforcement learning.

References

1. Abbeel, P., Coates, A., Ng, A.Y.: Autonomous helicopter aerobatics through apprenticeship learning. *Int. J. Robot. Res.* **29**(13), 1608–1639 (2010)
2. Albrecht, S., Ramirez-Amaro, K., Ruiz-Ugalde, F., Weikersdorfer, D., Leibold, M., Ulbrich, M., Beetz, M.: Imitating human reaching motions using physically inspired optimization principles. In: *Proceedings of HUMANOIDS* (2011)
3. Argall, B.D., Chernova, S., Veloso, M., Browning, B.: A survey of robot learning from demonstration. *Robot. Auton. Syst.* **57**(5), 469–483 (2009)
4. Calinon, S., Alizadeh, T., Caldwell, D.G.: On improving the extrapolation capability of task-parameterized movement models. In: *Proceedings of IROS* (2013)
5. Hansen, N., Ostermeier, A.: Completely derandomized self-adaptation in evolution strategies. *Evol. Comput.* **9**(2), 159–195 (2001)
6. Jetchev, N., Toussaint, M.: TRIC: task space retrieval using inverse optimal control. *Auton. Robot.* **37**(2), 169–189 (2014)
7. Kalakrishnan, M., Pastor, P., Righetti, L., Schaal, S.: Learning objective functions for manipulation. In: *Proceedings of ICRA* (2013)

8. Kolter, J.Z., Abbeel, P., and Ng, A.Y. Hierarchical apprenticeship learning with application to quadruped locomotion. In: NIPS (2008)
9. Kroemer, O., van Hoof, H., Neumann, G., Peters, J.: Learning to predict phases of manipulation tasks as hidden states. In: Proceedings of ICRA (2014)
10. Levine, S., Koltun, V.: Continuous inverse optimal control with locally optimal examples. In: Proceedings of ICML (2012)
11. Mombaur, K., Truong, A., Laumond, J.-P.: From human to humanoid locomotion an inverse optimal control approach. *Auton. Robot.* **28**(3), 369–383 (2010)
12. Muhlig, M., Gienger, M., Steil, J.J., Goerick, C.: Automatic selection of task spaces for imitation learning. In: Proceedings of IROS (2009)
13. Nocedal, J., Wright, S.J.: Numerical Optimization. Springer, New York (2006)
14. Paraschos, A., Daniel, C., Peters, J., Neumann, G.: Probabilistic Movement Primitives. In: NIPS (2013)
15. Pastor, P., Kalakrishnan, M., Chitta, S., Theodorou, E., Schaal, S.: Skill learning and task outcome prediction for manipulation. In: Proceedings of ICRA (2011)
16. Puydupin-Jamin, A.-S., Johnson, M., Bretl, T.: A convex approach to inverse optimal control and its application to modeling human locomotion. In: Proceedings of ICRA (2012)
17. Rückert, E.A., Neumann, G., Toussaint, M., Maass, W.: Learned graphical models for probabilistic planning provide a new class of movement primitives. *Front. Comput. Neurosci.* **6**, 1–20 (2013)
18. Russell, S.: Learning agents for uncertain environments. In: Proceedings of the Conference on Computational Learning Theory (1998)
19. Schaal, S., Ijspeert, A., Billard, A.: Computational approaches to motor learning by imitation. *Philos. Trans. R. Soc. Lond.* **358**, 537–547 (2003)
20. Stulp, F., Raiola, G., Hoarau, A., Ivaldi, S., Sigaud, O.: Learning compact parameterized skills with a single regression. In: Proceedings of HUMANOIDS (2013)
21. Theodorou, E., Buchli, J., Schaal, S.: A generalized path integral control approach to reinforcement learning. *J. Mach. Learn. Res.* **11**, 3137–3181 (2010)
22. Toussaint, M.: Newton methods for k-order markov constrained motion problems (2014). [arXiv:1407.0414](https://arxiv.org/abs/1407.0414) [cs.RO]
23. Toussaint, M., Ratliff, N., Bohg, J., Righetti, L., Englert, P., Schaal, S.: Dual execution of optimized contact interaction trajectories. In: Proceedings of IROS (2014)
24. Zhifei, S., Joo, E.M.: A survey of inverse reinforcement learning techniques. *Int. J. Intell. Comput. Cybern.* **5**(3), 293–311 (2012)
25. Ziebart, B.D., Maas, A., Bagnell, J.A., Dey, A.K.: Maximum entropy inverse reinforcement learning. In: Proceedings of the AAAI Conference on Artificial Intelligence (2008)