

Goal-Driven Dynamics Learning via Bayesian Optimization

Somil Bansal, Roberto Calandra, Ted Xiao, Sergey Levine, and Claire J. Tomlin

Abstract—Real-world robots are becoming increasingly complex and commonly act in poorly understood environments where it is extremely challenging to model or learn their true dynamics. Therefore, it might be desirable to take a task-specific approach, wherein the focus is on explicitly learning the dynamics model which achieves the best control performance for the task at hand, rather than learning the true dynamics. In this work, we use Bayesian optimization in an active learning framework where a locally linear dynamics model is learned with the intent of maximizing the control performance, and used in conjunction with optimal control schemes to efficiently design a controller for a given task. This model is updated directly based on the performance observed in experiments on the physical system in an iterative manner until a desired performance is achieved. We demonstrate the efficacy of the proposed approach through simulations and real experiments on a quadrotor testbed.

I. INTRODUCTION

Given the system dynamics, optimal control schemes such as LQR, MPC, and feedback linearization can efficiently design a controller that maximizes a performance criterion. However, depending on the system complexity, it can be quite challenging to model its true dynamics. Moreover, for a given task, a globally accurate dynamics model is not always necessary to design a controller. Often, partial knowledge of the dynamics is sufficient, e.g., for trajectory tracking purposes a local linearization of a non-linear system is often sufficient. In this paper we argue that, for complex systems, it might be preferable to adapt the controller design process for the specific task, using a learned system dynamics model sufficient to achieve the desired performance.

We propose Dynamics Optimization via Bayesian Optimization (aDOBO), a Bayesian Optimization (BO) based active learning framework to learn the dynamics model that achieves the best performance for a given task based on the performance observed in experiments on the physical system. This active learning framework takes into account all past experiments and suggests the next experiment in order to learn the most about the relationship between the performance criterion and the model parameters. Particularly important for robotic systems is the use of data-efficient approaches, where only few experiments are needed to obtain improved performance. Hence, we employ BO, an optimization method often used to optimize a performance criterion

All authors are with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley. {somil, roberto.calandra, t.xiao, svlevine, tomlin}@eecs.berkeley.edu

*This research is supported by NSF under the CPS Frontiers VehiCal project (1545126), by the UC-Philippine-California Advanced Research Institute under project IIID-2016-005, and by the ONR MURI Embedded Humans (N00014-16-1-2206).

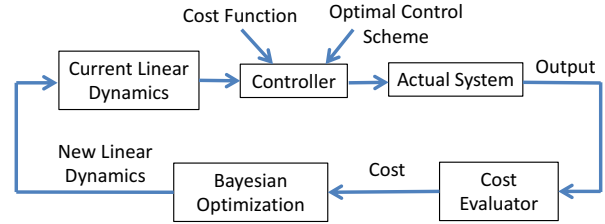


Fig. 1: aDOBO: A Bayesian optimization-based active learning framework for optimizing the dynamics model for a given cost function, directly based on the observed cost values.

while keeping the number of evaluations of the physical system small [1]. Specifically, we use BO to optimize the dynamics model with respect to the desired task, where the dynamics model is updated after every experiment so as to maximize the performance on the physical system. A flow diagram of our framework is shown in Figure 1. The current linear dynamics model, together with the cost function (also referred to as task/performance criterion), are used to design a controller with an appropriate optimal control scheme. The cost (or performance) of the controller is evaluated in closed-loop operation with the actual (unknown) physical plant. BO uses this performance information to iteratively update the dynamics model to improve the performance. This procedure corresponds to optimizing the linear system dynamics with the purpose of maximizing the performance of the final controller. Hence, unlike traditional system identification approaches, it does not necessarily correspond to finding the most accurate dynamics model, but rather the model yielding the best controller performance when provided to the optimal control method used.

Traditional system identification approaches are divided into two stages: 1) creating a dynamics model by minimizing some prediction error (e.g., using least squares) 2) using this dynamics model to generate an appropriate controller. In this approach, modeling the dynamics can be considered an offline process as there is no information flow between the two design stages. In online methods, the dynamics model is instead iteratively updated using the new data collected by evaluating the controller [2]. Our approach is an online method. Both for the online and the offline cases, creating a dynamics model based only on minimizing the prediction error can introduce sufficient inaccuracies to lead to suboptimal control performance [3]–[8]. Using machine learning techniques, such as Gaussian processes, does not alleviate this issue [9]. Instead, authors in [3] proposed to optimize the dynamics model directly with respect to the controller performance, but since the dynamics model is optimized offline, the resultant model is not necessarily

optimal for the actual system. We instead explicitly find the dynamics model that produces the best control performance for the *actual* system.

Previous studies addressed the problem of optimizing a controller using BO. In [10]–[12], the authors tuned the penalty matrices in an LQR problem for performance optimization. Although interesting results emerge from these studies, tuning penalty matrices may not achieve the desired performance when an accurate system model is not available. Our approach overcomes these challenges as it does not rely on an accurate system dynamics model. In [13], the authors directly learn the parameters of a linear feedback controller using BO. However, a typical controller might be non-linear and can contain hundreds of parameters; it is not feasible to optimize such high-dimensional controllers using BO [1]. Since aDOBO does not aim at directly learning the controller, it is agnostic to the dimensionality of the controller. It can leverage the low-dimensional structure of the dynamics to optimize the high-dimensional controllers. Moreover, it does not need to impose any structure on the controller and can easily design general non-linear controllers as well (see Sec. IV). The problem of updating a system model to improve control performance is also related to adaptive control, where the model parameters are identified from sensor data, and subsequently the updated model is used to design a controller (see [14]–[18]). However, in adaptive control, the model parameters are generally updated to get a good prediction model and not necessarily to maximize the controller performance. In contrast, we explicitly take into account the observed performance and search for the model that achieves the highest performance.

To the best of our knowledge, this is the first method that optimizes a dynamics model to maximize the control performance on the actual system. Our approach does not require the prior knowledge of an accurate dynamics model, nor of its parameterized form. Instead, the dynamics model is optimized, in an active learning setting, directly with respect to the desired cost function using data-efficient BO.

II. PROBLEM FORMULATION

Consider an unknown, stable, discrete-time, potentially non-linear, dynamical system

$$z_{k+1} = f(z_k, u_k), \quad k \in \{0, 1, \dots, N-1\}, \quad (1)$$

where $z_k \in \mathbb{R}^{n_x}$ and $u_k \in \mathbb{R}^{n_u}$ denote the system state and the control input at time k respectively. Given an initial state z_0 , the objective is to design a controller that minimizes the cost function J subject to the dynamics in (1)

$$J_0^* = \min_{\mathbf{u}_0^{N-1}} J_0(\mathbf{z}_0^N, \mathbf{u}_0^{N-1}) = \min_{\mathbf{u}_0^{N-1}} \sum_{i=0}^{N-1} l(z_i, u_i) + g(z_N, u_N),$$

subject to $z_{k+1} = f(z_k, u_k)$,

(2)

where $\mathbf{z}_i^N := (z_i, z_{i+1}, \dots, z_N)$. \mathbf{u}_i^{N-1} is similarly defined. One of the key challenges in designing such a controller is the modeling of the unknown system dynamics in (1). In this work, we model (1) as a linear time-invariant (LTI) system

with system matrices (A_θ, B_θ) . The system matrices are parameterized by $\theta \in \mathcal{M} \subseteq \mathbb{R}^d$, which is to be varied during the learning procedure. For a given θ and the current system state z_k , let $\pi_k(z_k, \theta)$ denote the optimal control sequence for the *linear* system (A_θ, B_θ) for the horizon $\{k, k+1, \dots, N\}$

$$\pi_k(z_k, \theta) := \bar{\mathbf{u}}_k^{N-1} = \arg \min_{\mathbf{u}_k^{N-1}} J_k(\mathbf{z}_k^N, \mathbf{u}_k^{N-1}), \quad (3)$$

subject to $z_{j+1} = A_\theta z_j + B_\theta u_j$.

The key difference between (2) and (3) is that the controller is designed for the parameterized linear system as opposed to the true system. As θ is varied, different matrix pairs (A_θ, B_θ) are obtained, which result in different controllers $\pi(\cdot, \theta)$. Our aim is to find, among all linear models, the linear model $(A_{\theta^*}, B_{\theta^*})$ whose controller $\pi(\cdot, \theta^*)$ minimizes J_0 (ideally achieves J_0^*) for the *actual system*, i.e.,

$$\theta^* = \arg \min_{\theta \in \mathcal{M}} J_0(\mathbf{z}_0^N, \mathbf{u}_0^{N-1}), \quad (4)$$

subject to $z_{k+1} = f(z_k, u_k)$, $u_k = \pi_k^1(z_k, \theta)$,

where $\pi_k^1(z_k, \theta)$ denotes the 1st control in the sequence $\pi_k(z_k, \theta)$. To make the dependence on θ explicit, we refer to J_0 in (4) as $J(\theta)$ here on. Note that $(A_{\theta^*}, B_{\theta^*})$ in (4) may not correspond to an actual linearization of the system, but simply to the linear model that gives the best performance on the actual system when its optimal controller is applied in a *closed-loop* fashion on the actual physical plant.

We choose LTI modeling to reduce the number of parameters used to represent the system, and make the dynamics learning process data efficient. Linear modeling also allows to efficiently design the controller in (3) for general cost functions (e.g., using MPC for any convex cost J). In general, the effectiveness of linear modeling depends on both the system and the control objective. If f is linear, a linear model is trivially sufficient for any control objective. If f is non-linear, a linear model may not be sufficient for all control tasks; however, for regulation and trajectory tracking tasks, a linear model is often adequate (see Sec. V). A linear parameterization is also used in adaptive control for similar reasons [18]. Nevertheless, the proposed framework can handle more general model classes as long as the optimal control problem in (3) can be solved for those classes.

Since f is unknown, the shape of the cost function, $J(\theta)$, in (4) is unknown. The cost is thus evaluated empirically in each experiment, which is often expensive as it involves conducting an experiment. Thus, the goal is to solve the optimization problem in (4) with as few evaluations as possible. In this paper, we do so via BO.

III. BACKGROUND

In order to optimize (A_θ, B_θ) , we use BO. In this section, we briefly introduce Gaussian processes and BO.

A. Gaussian Process (GP)

Since the function $J(\theta)$ in (4) is unknown a priori, we use nonparametric GP models to approximate it over its domain \mathcal{M} . GPs are a popular choice for probabilistic non-parametric regression, where the goal is to find a nonlinear

map, $J(\theta) : \mathcal{M} \rightarrow \mathbb{R}$, from an input vector $\theta \in \mathcal{M}$ to the function value $J(\theta)$. Hence, we assume that function values $J(\theta)$, associated with different values of θ , are random variables and that any finite number of these random variables have a joint Gaussian distribution dependent on the values of θ [19]. For GPs, we define a prior mean function and a covariance function (or kernel), $k(\theta_i, \theta_j)$, which defines the covariance between any two function values, $J(\theta_i)$ and $J(\theta_j)$. In this work, the mean is assumed to be zero without loss of generality. The choice of kernel is problem-dependent and encodes general assumptions such as smoothness of the unknown function. In the experimental section, we employ the 5/2 Matérn kernel where the hyperparameters are optimized by maximizing the marginal likelihood [19]. This kernel function implies that the underlying function J is differentiable and takes values within the $2\sigma_f$ confidence interval with high probability.

The GP framework can be used to predict the distribution of the performance function $J(\theta^*)$ at an arbitrary input θ^* based on the past observations, $\mathcal{D} = \{\theta_i, J(\theta_i)\}_{i=1}^n$. Conditioned on \mathcal{D} , the mean and variance of the prediction are

$$\mu(\theta^*) = \mathbf{k}\mathbf{K}^{-1}\mathbf{J}; \quad \sigma^2(\theta^*) = k(\theta^*, \theta^*) - \mathbf{k}\mathbf{K}^{-1}\mathbf{k}^T, \quad (5)$$

where \mathbf{K} is the kernel matrix with $K_{ij} = k(\theta_i, \theta_j)$, $\mathbf{k} = [k(\theta_1, \theta^*), \dots, k(\theta_n, \theta^*)]$ and $\mathbf{J} = [J(\theta_1), \dots, J(\theta_n)]$. Thus, the GP provides both the expected value of the performance function at any arbitrary point θ^* as well as a notion of the uncertainty of this estimate.

B. Bayesian Optimization (BO)

Bayesian optimization aims to find the global minimum of an unknown function [1], [20], [21]. BO is particularly suitable for the scenarios where evaluating the unknown function is expensive, which fits our problem in Sec. II. At each iteration, BO uses the past observations \mathcal{D} to model the objective function, and uses this model to determine informative sample locations. A common model used in BO for the underlying objective, and the one that we consider, are Gaussian processes (see Sec. III-A). Using the mean and variance predictions of the GP from (5), BO computes the next sample location by optimizing the so-called acquisition function, $\alpha(\cdot)$. Different acquisition functions are used in literature to trade off between exploration and exploitation during the optimization process [1]. For example, the next evaluation for expected improvement (EI) acquisition function [22] is given by $\theta^* = \arg \min_{\theta} \alpha(\theta)$ where

$$\alpha(\theta) = \sigma(\theta)[u\Phi(u) + \phi(u)]; \quad u = (\mu(\theta) - T)/\sigma(\theta). \quad (6)$$

$\Phi(\cdot)$ and $\phi(\cdot)$ in (6), respectively, are the standard normal cumulative distribution and probability density functions. The target value T is the minimum of all explored data. Intuitively, EI selects the next parameter point where the expected improvement over T is maximal. Repeatedly evaluating the system at points given by (6) thus improves the observed performance. Note that optimizing $\alpha(\theta)$ in (6) does not require physical interactions with the system, but only evaluation of the GP model. When a new set of optimal

Algorithm 1: aDOBO algorithm

```

1  $\mathcal{D} \leftarrow$  if available:  $\{\theta, J(\theta)\}$ 
2 Prior  $\leftarrow$  if available: Prior of the GP hyperparameters
3 Initialize GP with  $\mathcal{D}$ 
4 while optimize do
5   Find  $\theta^* = \arg \min_{\theta} \alpha(\theta)$ ;  $\theta' \leftarrow \theta^*$ 
6    $\mathbf{z}_0^N = \{\}$ ,  $\mathbf{u}_0^{N-1} = \{\}$ 
7   for  $i = 0 : N - 1$  do
8     Given  $z_i$  and  $(A_{\theta'}, B_{\theta'})$ , compute  $\pi_i(z_i, \theta')$ 
9     Apply  $\pi_i^1(z_i, \theta')$  on the real system and
       measure  $z_{i+1}$ 
10     $\mathbf{z}_0^N \leftarrow (\mathbf{z}_0^N, z_{i+1})$ 
11     $\mathbf{u}_0^{N-1} \leftarrow (\mathbf{u}_0^{N-1}, \pi_i^1(z_i, \theta'))$ 
12  Evaluate  $J(\theta') := J_0(\mathbf{z}_0^N, \mathbf{u}_0^{N-1})$  using (2)
13  Update GP and  $\mathcal{D}$  with  $\{\theta', J(\theta')\}$ 

```

parameters θ^* is determined, they are finally evaluated on the real objective function J (i.e., the system).

IV. DYNAMICS OPTIMIZATION VIA BO (aDOBO)

This section presents the technical details of aDOBO, a novel framework for optimizing dynamics model for maximizing the resultant controller performance. In this work, $\theta \in \mathbb{R}^{n_x(n_x+n_u)}$, i.e., each dimension in θ corresponds to an entry of the A_{θ} or B_{θ} matrices. This parameterization is chosen for simplicity, but other parameterizations can easily be used.

Given an initial state of the system z_0 and the current system dynamics model $(A_{\theta'}, B_{\theta'})$, we design an optimal control sequence $\pi_0(z_0, \theta')$ that minimizes the cost function $J_0(\mathbf{z}_0^N, \mathbf{u}_0^{N-1})$, i.e., we solve the optimal control problem in (3). The first control of this control sequence is applied on the *actual system* and the next state z_1 is measured. We then similarly compute $\pi_1(z_1, \theta')$ starting at z_1 , apply the first control in the obtained control sequence, measure z_2 , and so on until we get z_N . Once \mathbf{z}_0^N and \mathbf{u}_0^{N-1} are obtained, we compute the true performance of \mathbf{u}_0^{N-1} on the actual system by analytically computing $J_0(\mathbf{z}_0^N, \mathbf{u}_0^{N-1})$ using (2). We denote this cost by $J(\theta')$ for simplicity. We next update the GP based on the collected data sample $\{\theta', J(\theta')\}$. Finally, we compute θ^* that minimizes the corresponding acquisition function $\alpha(\theta)$ and repeat the process for $(A_{\theta^*}, B_{\theta^*})$. Our approach is illustrated in Figure 1 and summarized in Algorithm 1. Intuitively, aDOBO directly learns the shape of the cost function $J(\theta)$ as a function of linearizations (A_{θ}, B_{θ}) . Instead of learning the global shape of this function through random queries, it analyzes the performance of all the past evaluations and by optimizing the acquisition function, generates the next query that provides the maximum information about the minima of the cost function. This direct *minima-seeking* behavior based on the *actual observed performance* ensures that our approach is data-efficient. Thus, in the space of all linearizations, we efficiently and directly search for the linearization whose corresponding controller minimizes J_0 on the actual system.

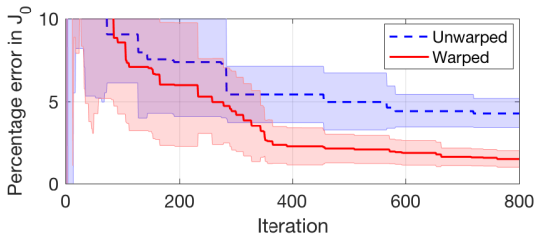


Fig. 2: Dubins car: mean and standard deviation of η during the learning process (over 10 trials). aDOBO reaches within the 10% of the optimal cost in just 100 iterations, starting from a random dynamics model. Using a log warping on the cost function further accelerates the learning.

Since the problem in (3) is an optimal control problem for the linear system $(A_{\theta'}, B_{\theta'})$, depending on the form of the cost function J , different optimal control schemes can be used. For example, if J is quadratic, the optimal controller is a linear feedback controller given by the solution of a Riccati equation. If J is a general convex function, the optimal control problem is solved through a general convex MPC solver, and the resultant controller could be non-linear. Thus, depending on the form of J , the controller designed by aDOBO can be linear or non-linear. This property causes aDOBO to perform well in the scenarios where a linear controller is not sufficient, as shown in Sec. VI-B. More generally, the proposed framework is modular and other control schemes can be used that are more suitable for a given cost function, which allows us to capture a richer controller space.

Note that the GP in our algorithm can be initialized with dynamics models whose controllers are known to perform well on the actual system. This generally leads to a faster convergence. For example, when a good linearization of the system is known, it can be used to initialize \mathcal{D} . When no information is known about the system a priori, the initial models are queried randomly. Finally, note that aDOBO can also be used when the real system is stochastic. In this case, aDOBO will minimize the expected cost.

V. NUMERICAL SIMULATIONS

In this section, we present some simulation results on the performance of the proposed method for controller design.

A. Dubins Car System

For the first simulation, we consider a three dimensional non-linear Dubins car whose dynamics are given as

$$\dot{x} = v \cos \phi, \quad \dot{y} = v \sin \phi, \quad \dot{\phi} = \omega, \quad (7)$$

where $z := (x, y, \phi)$ is the state of system, $p = (x, y)$ is the position, ϕ is the heading, v is the speed, and ω is the turn rate. The input (control) to the system is $u := (v, \omega)$. For simulation purposes, we discretize the dynamics at a frequency of 10Hz. Our goal is to design a controller that steers the system to the equilibrium point $z^* = 0, u^* = 0$ starting from the state $z_0 := (1.5, 1, \pi/2)$. In particular, we

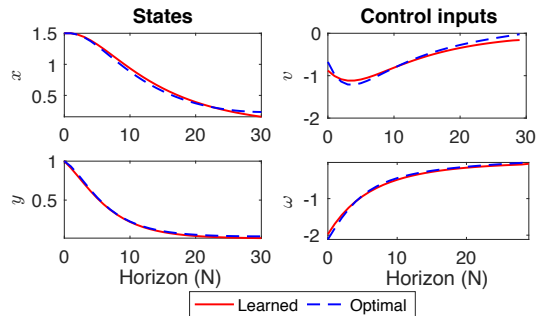


Fig. 3: Dubins car: state and control trajectories for the learned and the true system. The two trajectories are very similar, indicating that the learned dynamics model represents the system behavior accurately around the desired state.

want to minimize the cost function

$$J_0(z_0^N, u_0^{N-1}) = \sum_{k=0}^{N-1} (z_k^T Q z_k + u_k^T R u_k) + z_N^T Q_f z_N. \quad (8)$$

We choose $N = 30$. Q , Q_f and R are all chosen as identity matrices of appropriate sizes. We also assume that the dynamics are not known; hence, we cannot directly design a controller to steer the system to the desired equilibrium. Instead, we use aDOBO to find a linearization of dynamics in (7) that minimizes the cost function in (8), directly from the experimental data. In particular, we represent the system in (7) by a parameterized linear system $z_{k+1} = A_{\theta} z_k + B_{\theta} u_k$, design a controller for this system and apply it on the actual system. Based on the observed performance, BO suggests a new linearization and the process is repeated. Since the cost function is quadratic in this case, the optimal control problem for a particular θ is an LQR problem, and can be solved efficiently. For BO, we use the MATLAB library *BayesOpt* [23]. Since there are 3 states and 2 inputs, we learn 15 parameters in total, one corresponding to each entry of the A_{θ} and B_{θ} matrices. The bounds on the parameters are chosen randomly as $\mathcal{M} = [-2, 2]^{15}$. As acquisition function, we use EI (see eq. (6)). Since no information is assumed to be known about the system, the GP was initialized with a random θ . We also warp the cost function J using the *log* function before passing it to BO. Warping makes the cost function smoother while maintaining its monotonic properties, which makes the sampling process in BO more efficient and leads to a faster convergence.

For comparison, we compute the true optimal controller that minimizes (8) subject to the dynamics in (7) using the non-linear solver *fmincon* in MATLAB to get the minimum achievable cost J_0^* across all controllers. We use the percentage error between the true optimal cost J_0^* and the cost achieved by aDOBO as our comparison metric in this work

$$\eta_n = 100 \times (J_0^* - J(\theta_n)) / J_0^*, \quad (9)$$

where $J(\theta_n)$ is the best cost achieved by aDOBO by iteration n . In Fig. 2, we plot η_n for Dubins car. As learning progresses, aDOBO gathers more and more information about the minimum of J_0 and reaches within 10% of J_0^*

in just 100 iterations, demonstrating its effectiveness in designing a controller for an unknown system just from the experimental data. Fig. 2 also highlights the effect of warping in BO. A well warped function converges faster to the optimal performance. We also compared the control and state trajectories obtained from the learned controller with the optimal control and state trajectories. As shown in Fig. 3, the learned system matrices not only achieve the optimal cost, but also follow the optimal state and control trajectories very closely. Even though the trajectories are very close to each other for the true system and its learned linearization, this linearization may not correspond to any *actual* linearization of the system. The next simulation illustrates this key property of aDOBO more clearly.

B. A Simple 1D Linear System

For this simulation, we consider a simple 1D linear system

$$z_{k+1} = z_k + u_k, \quad (10)$$

where z_k and u_k are the state and the input of the system at time k . Although the dynamics model is very simple, it illustrates some key insights about the proposed method. Our goal is to design a controller that minimizes (8) starting from the state $z_0 = 1$. We choose $N = 30$ and $R = Q = Q_f = 1$. We assume that the dynamics are unknown and use aDOBO to learn the dynamics, where $\theta := (\theta_1, \theta_2) \in \mathbb{R}^2$ are the dynamics parameters to be learned.

The learning process converges in 45 iterations to the true optimal performance ($J_0^* = 1.61$), which is computed using LQR on the real system. The converged parameters are $\theta_1 = 1.69$ and $\theta_2 = 2.45$, which are vastly different from the true parameters $\theta_1 = 1$ and $\theta_2 = 1$, even though the actual system is a linear system. To understand this, we plot the cost obtained on the true system J_0 as a function of linearization parameters (θ_1, θ_2) in Fig. 4. Since the performances of the two sets of parameters are very close to each other, a direct performance based learning process (e.g., aDOBO) cannot distinguish between them and both sets are equally optimal for it. More generally, a wide range of parameters lead to similar performance on the actual system. Hence, we can expect the proposed approach to recover the optimal controller and the actual state trajectories, but not necessarily the true dynamics or its true linearization. This simulation also suggests that the true dynamics of the system may

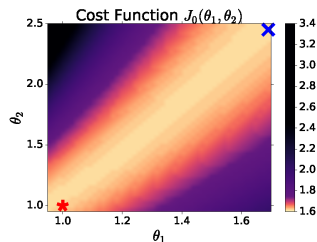


Fig. 4: Cost of the actual system in (10) as a function of the linearization parameters (θ_1, θ_2) . The parameters obtained by aDOBO (the pink X) yield to performance very close to the true system parameters (the green *). Note that aDOBO does not necessarily converge to the true parameters.

not even be required as far as the control performance is concerned.

C. Cart-pole System

We next apply aDOBO to a cart-pole system

$$\begin{aligned} (M + m)\ddot{x} - ml\ddot{\psi} \cos \psi + ml\dot{\psi}^2 \sin \psi &= F, \\ l\ddot{\psi} - g \sin \psi &= \ddot{x} \cos \psi, \end{aligned} \quad (11)$$

where x denotes the position of the cart with mass M , ψ denotes the pendulum angle, and F is a force that serves as the control input. The massless pendulum is of length l with a mass m attached at its end. Define the system state as $z := (x, \dot{x}, \psi, \dot{\psi})$ and the input as $u := F$. Starting from the state $(0, 0, \frac{\pi}{6}, 0)$, the goal is to keep the pendulum straight up, while keeping the state within given lower and upper bounds. In particular, we want to minimize the cost

$$\begin{aligned} J_0(z_0^N, u_0^{N-1}) &= \sum_{k=0}^{N-1} (z_k^T Q z_k + u_k^T R u_k) + z_N^T Q_f z_N \\ &+ \lambda \sum_{i=0}^N \max(0, \underline{z} - z_i, z_i - \bar{z}), \end{aligned} \quad (12)$$

where λ penalizes the deviation of state z_i below \underline{z} and above \bar{z} . We assume that the dynamics are unknown and use aDOBO to optimize the dynamics. For simulation, we discretize the dynamics at a frequency of 10Hz. We choose $N = 30$, $M = 1.5\text{Kg}$, $m = 0.175\text{Kg}$, $\lambda = 100$ and $l = 0.28\text{m}$. The $Q = Q_f = \text{diag}([0.1, 1, 100, 1])$ and $R = 0.1$ matrices are chosen to penalize the angular deviation significantly. We use $\underline{z} = [-2, -\infty, -0.1, -\infty]$ and $\bar{z} = [2, \infty, \infty, \infty]$, i.e., we are interested in controlling the pendulum while keeping the cart position within $[-2, 2]$, and limiting the pendulum overshoot to 0.1. The optimal control problem for a particular linearization is a convex MPC problem and solved using YALMIP [24]. The true J_0^* is computed using *fmincon*.

As shown in Fig. 5, aDOBO reaches within 20% of the optimal performance in 250 iterations and continue to make progress towards finding the optimal controller. This simulation demonstrates that the proposed method (a) is also applicable to highly non-linear systems, (b) can handle general convex cost functions that are not necessarily quadratic, and (c) different optimal control schemes can be used within the proposed framework. Since an MPC controller can in general be non-linear, this implies that the proposed method can also design complex non-linear controllers with an LTI parametrization.

VI. COMPARISON WITH OTHER METHODS

In this section, we compare our approach with some other online learning schemes for controller design.

A. Tuning (Q, R) vs aDOBO

In this section, we consider the case in which the cost function J_0 is quadratic (see Eq. (8)). Suppose that the actual linearization of the system around $z^* = 0$ and $u^* = 0$ is

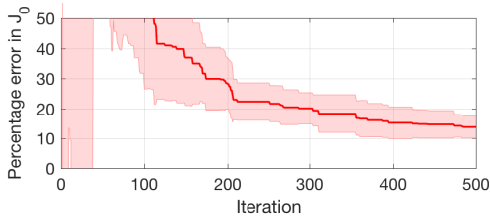


Fig. 5: Cart-pole system: mean and standard deviation of η during the learning process. The learned controller reaches within 20% of the optimal cost in 250 iterations, demonstrating the applicability of aDOBO to highly non-linear systems.

known and given by (A^*, B^*) . To design a controller for the actual system in such a case, it is a common practice to use an LQR controller for the linearized dynamics. However, the resultant controller may be sub-optimal for the actual non-linear system. To overcome this problem, authors in [10], [11] propose to optimize the controller by tuning penalty matrices Q and R in (8). In particular, we solve

$$\begin{aligned} \theta^* &= \arg \min_{\theta \in \mathcal{M}} J_0(\mathbf{z}_0^N, \mathbf{u}_0^{N-1}), \\ \text{sub. to } z_{k+1} &= f(z_k, u_k), \quad u_k = K(\theta)z_k, \\ K(\theta) &= LQR(A^*, B^*, W_Q(\theta), W_R(\theta), Q_f), \end{aligned} \quad (13)$$

where $K(\theta)$ denotes the LQR feedback matrix obtained for the system matrices (A^*, B^*) with W_Q and W_R as state and input penalty matrices, and can be computed analytically. For further details of LQR method, we refer interested readers to [25]. The difference between optimization problems (4) and (13) is that now we parameterize penalty matrices W_Q and W_R instead of system dynamics. The optimization problem in (13) is solved using BO in a similar fashion as we solve (4) [10]. The parameter θ , in this case, can be initialized by the actual penalty matrices Q and R , instead of a random query, which generally leads to a much faster convergence. An alternative approach is to use aDOBO, except that now we can use (A^*, B^*) as initializations for the system matrices A and B . Actual penalty matrices Q and R are used for aDOBO.

When (A^*, B^*) are known to a good accuracy, (Q, R) tuning method is expected to converge quickly to the optimal performance compared to aDOBO as it needs to learn fewer parameters, i.e., $(n_x + n_u)$ (assuming diagonal penalty matrices) compared to $n_x(n_x + n_u)$ parameters for aDOBO. However, when there is error in (A^*, B^*) (or more generally if dynamics are unknown), the performance of the (Q, R) tuning method can degrade significantly as it relies on an accurate linearization of the system dynamics, rendering the method impractical for control design purposes. To compare the two methods we use the Dubins car model in Eq. (7). The rest of the simulation parameters are same as Section V-A. We compute the linearization of Dubins car around $z^* = 0$ and $u^* = 0$ using (7) and add random matrices (A_r, B_r) to them to generate $A' = (1 - \alpha)A^* + \alpha A_r$ and $B' = (1 - \alpha)B^* + \alpha B_r$. We then initialize both methods with (A', B') for different α s. As shown in Fig. 6, the (Q, R) tuning method outperforms aDOBO, when there is no noise

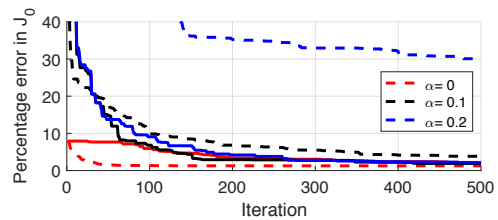


Fig. 6: Dubins car: Comparison between tuning the penalty matrices (Q, R) [10] (dashed curves), and aDOBO (solid curves) for different noise levels in (A^*, B^*) , the true linearized dynamics around the desired goal state. When the true linearized dynamics are known perfectly, the (Q, R) tuning method outperforms aDOBO because fewer parameters are to be learned. Its performance, however, drops significantly as noise increases, rendering the method impractical for the scenarios where system dynamics are not known to a good accuracy.

in (A^*, B^*) . But as α increases, its performance deteriorates significantly. In contrast, aDOBO is fairly indifferent to the noise level, as it does not assume any prior knowledge of system dynamics. The only information assumed to be known is penalty matrices (Q, R) , which are generally designed by the user and hence are known a priori.

B. Learning K vs aDOBO

When the cost function is quadratic, another potential approach is to directly parameterize and optimize the feedback matrix $K \in \mathbb{R}^{n_x n_u}$ in (13) [13] as

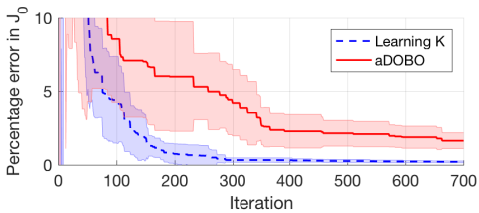
$$\begin{aligned} \theta^* &= \arg \min_{\theta \in \mathcal{M}} J_0(\mathbf{z}_0^N, \mathbf{u}_0^{N-1}), \\ \text{sub. to } z_{k+1} &= f(z_k, u_k), \quad u_k = K_\theta z_k. \end{aligned} \quad (14)$$

The advantage of this approach is that only $n_x n_u$ parameters are learned compared to $n_x(n_x + n_u)$ parameters in aDOBO, which is also evident from Fig. 7a, wherein the learning process for K converges much faster than that for aDOBO. However, a linear controller might not be sufficient for general cost functions, and non-linear controllers are required to achieve a desired performance. As shown in Sec. V-C, aDOBO is not limited to linear controllers; hence, it outperforms the K learning method in such scenarios. Consider, for example, the linear system

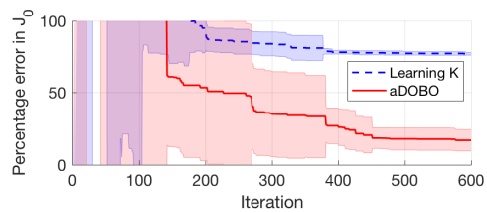
$$x_{k+1} = x_k + y_k, \quad y_{k+1} = y_k + u_k, \quad (15)$$

and the cost function in Eq. (12) with state $z_k = (x_k, y_k)$, $N = 30$, $\underline{z} = [0.5, -0.4]$ and $\bar{z} = [\infty, \infty]$. Q , Q_f and R are all identity matrices of appropriate sizes, and $\lambda = 100$.

As evident from Fig. 7b, directly learning a feedback matrix performs poorly with an error as high as 80% from the optimal cost. Since the cost is not quadratic, the optimal controller is not necessarily linear; however, since the controller in (14) is restricted to a linear space, it performs rather poorly in this case. In contrast, aDOBO continues to improve performance and reaches within 20% of the optimal cost within few iterations, because we implicitly parameterize a much richer controller space via learning A



(a) Dubins car



(b) System of Eq. (15)

Fig. 7: Mean and standard deviation of η obtained via directly learning the feedback controller K [13] and aDOBO for different cost functions. (a) Comparison for the quadratic cost function of Eq. (8). Directly learning K converges to the optimal performance faster because fewer parameters are to be learned. (b) Comparison for the non-quadratic cost function of Eq. (12). Since the optimal controller for the actual system is not necessarily linear in this case, directly learning K leads to a poor performance

Iteration	aDOBO	Learning Control Sequence
200	53 \pm 50 %	605 \pm 420%
400	27 \pm 12 %	357 \pm 159%
600	17 \pm 7 %	263 \pm 150%

TABLE I: System in (15): mean and standard deviation of η for aDOBO, and for directly learning the control sequence. Since the space of control sequence is huge, the error is substantial even after 600 iterations.

and B . In this example, we capture non-linear controllers by using a linear dynamics model with a convex MPC solver. Since the underlying system is linear, the true optimal controller is also in our search space. Our algorithm makes sure that we make a steady progress towards finding that controller. However, we are not restricted to learning a linear controller K . One can also directly learn the actual control sequence to be applied to the system (which also captures the optimal controller). This approach may not be data-efficient compared to aDOBO as the control sequence space can be very large depending on the problem horizon, and will require a large number of experiments. As shown in Table I, the performance error is more than 250% even after 600 iterations, rendering the method impractical for real systems.

C. Adaptive Control vs aDOBO

In this work, we aim to directly find the best linearization based on the observed performance. Another approach is to learn a true linearization of the system based on the observed state and input trajectory during the experiments. The underlying hypothesis is that as more and more data is collected, a better linearization is obtained, eventually leading to an improved control performance. This approach is in-line with the traditional model identification and the adaptive control frameworks. Let $({}_j\mathbf{z}_0^N, {}_j\mathbf{u}_0^{N-1})$ denotes the state and input trajectories for experiment j . We also let $\mathcal{D}_i = \cup_{j=1}^i ({}_j\mathbf{z}_0^N, {}_j\mathbf{u}_0^{N-1})$. After experiment i , we fit an LTI model of the form $z_{k+1} = A_i z_k + B_i u_k$ using least squares on data in \mathcal{D}_i and then use this model to obtain a new controller for experiment $i+1$. We apply the approach on the linear system in (15) and the non-linear system in (7) with the cost function in (8). For the linear system, the approach converges to the true system dynamics in 5 iterations. However, this approach

Iteration	aDOBO	Learning via LS
200	6 \pm 3.7 %	166.7 \pm 411%
400	2.2 \pm 1.1 %	75.9 \pm 189%
600	1.8 \pm 0.7 %	70.7 \pm 166%

TABLE II: Dubins car: mean and standard deviation of η obtained via learning (A, B) through least squares (LS), and through aDOBO.

performs rather poorly on the non-linear system, as shown in Table II. When the underlying system is non-linear, all state and input trajectories may not contribute to the performance improvement. A good linearization should be obtained from the state and input trajectories in the region of interest, which depends on the task. For example, if we want to regulate the system to the equilibrium $(0, 0)$, a linearization of the system around $(0, 0)$ should be obtained. Thus, it is desirable to use the system trajectories that are close to this equilibrium point. However, a naive prediction error based approach has no means to select these “good” trajectories from the pool of trajectories and hence can lead to a poor performance. In contrast, aDOBO does not suffer from these limitations, as it explicitly utilizes a performance based optimization. A summary of the advantages and limitations of the four methods is provided in Table III.

VII. QUADROTOR POSITION TRACKING EXPERIMENTS

We now present the results of our experiments on Crazyflie 2.0, which is an open source nano quadrotor platform developed by Bitcraze. Its small size, low cost, and robustness make it an ideal platform for testing

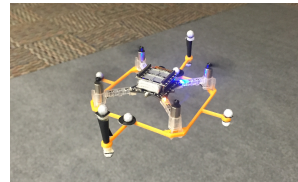


Fig. 8: The Crazyflie 2.0

new control paradigms. Recently, it has been extensively used to demonstrate aggressive flights [26], [27]. For small yaw, the quadrotor system is modeled as a rigid body with a ten dimensional state vector $s := [p, v, \zeta, \omega]$, which includes the position $p = (x, y, z)$ in an inertial frame I , linear velocities $v = (v_x, v_y, v_z)$ in I , attitude (orientation) represented by Euler angles ζ , and angular velocities ω . The system is controlled via three inputs $u := [u_1, u_2, u_3]$, where u_1 is the thrust along the z -axis, and u_2 and u_3 are

Method	Advantages	Limitations
(Q, R) learning [10]	Only $(n_x + n_u)$ parameters are to be learned so learning will be faster.	Performance can degrade significantly if the dynamics are not known to a good accuracy.
F learning [13]	Only $n_x n_u$ parameters are to be learned so learning will be faster.	Approach may not perform well for non-quadratic cost functions.
(A, B) learning via least squares	Can lead to a faster convergence when the underlying system is linear	Approach is not suitable for non-linear system.
aDOBO	Does not require any prior knowledge of system dynamics. Applicable to general cost functions.	Number of parameters to be learned is higher, i.e., $(n_x^2 + n_x n_u)$.

TABLE III: Relative advantages and limitations of different methods for automatic controller design.

rolling, pitching moments respectively. The full non-linear dynamics of a quadrotor are derived in [28], and its physical parameters are computed in [26]. Our goal in this experiment is to track a desired position p^* starting from the initial position $p_0 = [0, 0, 1]$. Formally, we minimize

$$J_0(\bar{s}_0^N, \mathbf{u}_0^{N-1}) = \sum_{k=0}^{N-1} (\bar{s}_k^T Q \bar{s}_k + u_k^T R u_k) + \bar{s}_N^T Q_f \bar{s}_N, \quad (16)$$

where $\bar{s} := [p - p^*, v, \zeta, \omega]$. Given the dynamics in [28], the desired optimal control problem can be solved using LQR; however, the resultant controller may still not be optimal for the actual system because (a) true underlying system is non-linear and (b) the actual system may not follow the dynamics in [28] exactly due to several unmodeled effects, as illustrated in our results. Hence, we assume that the dynamics of v_x and v_y are unknown, and model them as

$$\begin{bmatrix} f_{v_x} \\ f_{v_y} \end{bmatrix} = A_\theta \begin{bmatrix} \phi \\ \psi \end{bmatrix} + B_\theta u_1, \quad (17)$$

where A and B are parameterized through θ . Our goal is to learn the parameter θ^* that minimizes the cost in (16) for the *actual* Crazyflie using aDOBO. We use $N = 400$; the penalty matrix Q is chosen to penalize the position deviation. In our experiments, Crazyflie was flown in presence of a VICON motion capture system, which along with on-board sensors provides the full state information at 100Hz. The optimal control problem for a particular linearization in (17) is solved using LQR. For comparison, we compute the *nominal* optimal controller using the full dynamics in [28]. Figure 9 shows the performance of the controller from aDOBO compared with the nominal controller during the learning process. The nominal controller outperforms the learned controller initially, but within a few iterations, aDOBO performs better than the controller derived from the known dynamics model of Crazyflie. This is because aDOBO optimizes controller based on the performance of the *actual* system and hence can account for unmodeled effects. In 45 iterations, the learned controller outperforms the nominal controller by 12%, demonstrating the performance potential of aDOBO on real systems.

VIII. CONCLUSIONS AND FUTURE WORK

In this work, we introduce aDOBO, an active learning framework to optimize the system dynamics with the intent of maximizing the controller performance. Through simulations and real-world experiments, we demonstrate that

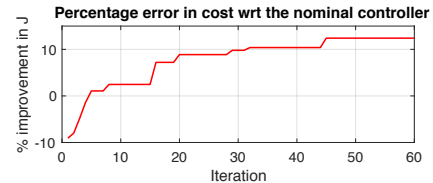


Fig. 9: Crazyflie: percentage error between the learned and the nominal controller. The nominal controller is obtained by using the full 12D non-linear dynamics model of the quadrotor. As learning progresses, aDOBO outperforms the nominal controller by 12% on the actual system, indicating the capability of aDOBO to overcome modeling errors.

aDOBO achieves the optimal control performance even when no prior information is known about the system dynamics.

Several interesting future directions emerge from this work. The dynamics model learned through aDOBO maximizes the performance on a single task. The obtained dynamics model thus may not necessarily perform well on a similar but different task. It will be interesting to generalize aDOBO to optimize the dynamics for a class of tasks, e.g., regulating to different states. Leveraging the state and input trajectory data, along with the observed performance, to further increase the data-efficiency of the learning process is another promising direction. During the learning process, aDOBO can query parameters which might lead to an unstable behavior on the actual system and can cause safety concerns. In such cases, it might be desirable to combine aDOBO with the exploration methods that explicitly take safety into account, e.g., such as SafeOpt [29], [30]. Finally, since BO is not scalable to higher-dimensional systems (roughly beyond 30-40 parameters) [1], exploring alternative ways to scale aDOBO to more complex non-linear dynamics models is another interesting direction.

REFERENCES

- [1] B. Shahriari, K. Swersky, Z. Wang, R. P. Adams, and N. de Freitas, "Taking the human out of the loop: A review of Bayesian optimization," *Proceedings of the IEEE*, vol. 104, no. 1, pp. 148–175, 2016.
- [2] M. P. Deisenroth, D. Fox, and C. E. Rasmussen, "Gaussian processes for data-efficient learning in robotics and control," *Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 2015.
- [3] J. Joseph, A. Geramifard, J. W. Roberts, J. P. How, and N. Roy, "Reinforcement learning with misspecified model classes," in *International Conference on Robotics and Automation*, 2013, pp. 939–946.
- [4] P. L. Donti, B. Amos, and J. Z. Kolter, "Task-based end-to-end model learning," *arXiv preprint arXiv:1703.04529*, 2017.
- [5] C. G. Atkeson, "Nonparametric model-based reinforcement learning," in *Advances in neural information processing systems*, 1998, pp. 1008–1014.

- [6] P. Abbeel, M. Quigley, and A. Y. Ng, "Using inaccurate models in reinforcement learning," in *International conference on Machine learning, 2006*, pp. 1–8.
- [7] M. Gevers, "Identification for control: From the early achievements to the revival of experiment design," *European journal of control*, vol. 11, no. 4-5, 2005.
- [8] H. Hjalmarsson, M. Gevers, and F. De Bruyne, "For model-based control design, closed-loop identification gives better performance," *Automatica*, vol. 32, no. 12, 1996.
- [9] D. Nguyen-Tuong and J. Peters, "Model learning for robot control: a survey," *Cognitive Processing*, vol. 12, no. 4, pp. 319–340, 2011.
- [10] A. Marco, P. Hennig, J. Bohg, S. Schaal, and S. Trimpe, "Automatic LQR tuning based on Gaussian process global optimization," in *International Conference on Robotics and Automation*, 2016.
- [11] S. Trimpe, A. Millane, S. Doessegger, and R. D'Andrea, "A self-tuning LQR approach demonstrated on an inverted pendulum," *IFAC Proceedings Volumes*, vol. 47, no. 3, pp. 11 281–11 287, 2014.
- [12] J. W. Roberts, I. R. Manchester, and R. Tedrake, "Feedback controller parameterizations for reinforcement learning," in *Symposium on Adaptive Dynamic Programming And Reinforcement Learning (ADPRL)*, 2011, pp. 310–317.
- [13] R. Calandra, A. Seyfarth, J. Peters, and M. P. Deisenroth, "Bayesian optimization for learning gaits under uncertainty," *Annals of Mathematics and Artificial Intelligence*, vol. 76, no. 1, pp. 5–23, 2015.
- [14] K. J. Åström and B. Wittenmark, *Adaptive control*. Courier Corporation, 2013.
- [15] M. Grimble, "Implicit and explicit LQG self-tuning controllers," *Automatica*, vol. 20, no. 5, pp. 661–669, 1984.
- [16] D. Clarke, P. Kanjilal, and C. Mohtadi, "A generalized LQG approach to self-tuning control part i. aspects of design," *International Journal of Control*, vol. 41, no. 6, pp. 1509–1523, 1985.
- [17] R. Murray-Smith and D. Sbarbaro, "Nonlinear adaptive control using nonparametric Gaussian process prior models," *IFAC Proceedings Volumes*, vol. 35, no. 1, pp. 325–330, 2002.
- [18] S. Sastry and M. Bodson, *Adaptive control: stability, convergence and robustness*. Courier Corporation, 2011.
- [19] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*. The MIT Press, 2006.
- [20] H. J. Kushner, "A new method of locating the maximum point of an arbitrary multipeak curve in the presence of noise," *Journal of Basic Engineering*, vol. 86, p. 97, 1964.
- [21] M. A. Osborne, R. Garnett, and S. J. Roberts, "Gaussian processes for global optimization," in *Learning and Intelligent Optimization (LION3)*, 2009, pp. 1–15.
- [22] J. Moćkus, "On bayesian methods for seeking the extremum," in *Optimization Techniques IFIP Technical Conference*, 1975.
- [23] R. Martinez-Cantin, "BayesOpt: a Bayesian optimization library for nonlinear optimization, experimental design and bandits," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 3735–3739, 2014.
- [24] J. Lofberg, "YALMIP: A toolbox for modeling and optimization in MATLAB," in *International Symposium on Computer Aided Control Systems Design*, 2005, pp. 284–289.
- [25] D. J. Bender and A. J. Laub, "The linear-quadratic optimal regulator for descriptor systems: discrete-time case," *Automatica*, 1987.
- [26] B. Landry, "Planning and control for quadrotor flight through cluttered environments," Master's thesis, MIT, 2015.
- [27] S. Bansal, A. K. Akametalu, F. J. Jiang, F. Laine, and C. J. Tomlin, "Learning quadrotor dynamics using neural network for flight control," in *Conference on Decision and Control*, 2016, pp. 4653–4660.
- [28] N. Abas, A. Legowo, and R. Akmeliawati, "Parameter identification of an autonomous quadrotor," in *International Conference On Mechatronics*, 2011, pp. 1–8.
- [29] Y. Sui, A. Gotovos, J. Burdick, and A. Krause, "Safe exploration for optimization with Gaussian processes," in *International Conference on Machine Learning*, 2015, pp. 997–1005.
- [30] F. Berkenkamp, A. P. Schoellig, and A. Krause, "Safe controller optimization for quadrotors with Gaussian processes," in *International Conference on Robotics and Automation*, 2016, pp. 491–496.