arXiv:1906.11003v2 [eess.SY] 31 Jul 2019

Chance-Constrained Trajectory Optimization for Non-linear Systems with Unknown Stochastic Dynamics

Onur Celik¹, Hany Abdulsamad² and Jan Peters^{2,3}

Abstract-Iterative trajectory optimization techniques for non-linear dynamical systems are among the most powerful and sample-efficient methods of model-based reinforcement learning and approximate optimal control. By leveraging time-variant local linear-quadratic approximations of system dynamics and reward, such methods can find both a target-optimal trajectory and time-variant optimal feedback controllers. However, the local linear-quadratic assumptions are a major source of optimization bias that leads to catastrophic greedy updates, raising the issue of proper regularization. Moreover, the approximate models' disregard for any physical state-action limits of the system causes further aggravation of the problem, as the optimization moves towards unreachable areas of the stateaction space. In this paper, we address the issue of constrained systems in the scenario of online-fitted stochastic linear dynamics. We propose modeling state and action physical limits as probabilistic chance constraints linear in both state and action and introduce a new trajectory optimization technique that integrates these probabilistic constraints by optimizing a relaxed quadratic program. Our empirical evaluations show a significant improvement in learning robustness, which enables our approach to perform more effective updates and avoid premature convergence observed in state-of-the-art algorithms.

I. INTRODUCTION

Model-based reinforcement learning has played an important role in the latest surge of popular research interest in learning-control of autonomous systems [1]. More specifically, trajectory-centric optimization techniques of non-linear dynamics have proven to be extremely sample efficient in comparison to model-free policy search approaches [2]–[4].

With the notable exception of [5], model-based trajectory optimization techniques [6], [7] are closely related to differential dynamic programming methods (DDP), initially presented in [8] and further generalized in [9]. DDP is a powerful approach for generating optimal trajectories with optimal time-variant feedback controllers. By relying on linear-quadratic approximations of the dynamics and reward around a nominal trajectory, DDP-based methods can leverage the local approximations to iteratively optimize both the trajectory and tracking feedback controllers in closed-form via dynamic programming [10]. This view of control has a computational advantage over direct optimization techniques such as collocation methods, which solve large optimization problems directly in the trajectory space and generally result only in open-loop control sequences [11].

However, despite the overwhelming success of DDP, it still suffers from multiple shortcomings. On the one hand, the greedy exploitation of poor local approximations of dynamics is a major problem that leads to premature convergence. This issue has been effectively addressed in recent research by proposing different schemes of regularization [2], [6], [7]. On the other hand, state and action constraints present a serious challenge, as they introduce hard non-linearities, that cannot be straightforwardly incorporated into the dynamic programming framework. The effect of constraints becomes more severe in settings where a global model is not available for automatic differentiation, hence requiring the linear approximation of the dynamics to be fitted online from samples.

We view these issues of DDP as interlocked. The inability of time-variant local linear models to consider state and action constraints results in updates that exploit unreachable parts of the state-actions space, leading to catastrophically poor linear-quadratic approximations in regions subject to hard non-linearities. Moreover, considering constraints becomes more challenging in scenarios with stochastic dynamics, in that the true state of the system is hidden and only available through sufficient statistics. Another crucial aspect in a stochastic setting is the infinite support of the noisy measurements, which results in the constraints being active over the whole state-action space.

To address these issues, we propose an augmented view of DDP that introduces the physical limits as probabilistic chance constraints linear in state and action. When considering time-variant linear-Gaussian approximations of the dynamic, we can relax the generally non-convex chance constraints by applying Boole's inequality. This relaxation allows us to formulate an additional quadratic program that forces the optimized nominal trajectory to stay in a feasible state-action region with high probability, all while considering the feedback gains optimized by DDP.

Several approaches to trajectory optimization for nonlinear systems address the problem of constrained dynamics on different levels. In the domain of deterministic environments, Tassa et al. considered action box-constraints in [12], while the authors in [13] introduce soft state-action limits via a Lagrange function augmentation. More sophisticated integration of constraints is presented in [14], in which the authors formulate a quadratic program to determine the active set of constraints at every iteration. In a stochastic setting, the work by Van Den Berg et. al [15] introduces probabilistic constraints as direct penalty terms on the cost function.

^{*}This work has received funding from the European Unions Horizon 2020 research and innovation program under grant agreement # 640554.

¹Onur Celik is with the Department of Computer Science, Universität Tübingen. mevluet-onur.celik@uni-tuebingen.de

²Hany Abdulsamad and Jan Peters are with the Department of Computer Science, Intelligent Autonomous Systems, Technische Universität Darmstadt. {abdulsamad, peters}@ias.tu-darmstadt.de

³Jan Peters is with the Max Planck Institute for Intelligent Systems.

Furthermore, probabilistic constraints are considered in the context of linear optimal control. In [16] the authors optimally handle probabilistic constraints by ellipsoidal relaxation for finite-horizon open-loop scenarios, while in [17] a similar problem is tackled by applying Boole's inequality. In [18] Vitus et al. propose an algorithm to extend the work in [17] and [16] by considering closed-loop uncertainty and optimizing the risk allocation. Finally, in [19] the problem of infeasible initial solutions is addressed by progressively introducing the constraint into the objective.

We situate our contribution in the class of differential dynamic programming for stochastic non-linear systems subject to probabilistic constraints in state and action. We empirically show that our proposed approach can deal with highly nonlinear constrained dynamic environments, leading to better overall performance and a robust learning process by virtue of improved online-fitted local approximations.

II. CHANCE-CONSTRAINED OPTIMIZATION

Chance constraints arise naturally in different fields of optimization when considering stochastic systems. For an overview, we refer to [20]. Dealing with such probabilistic constraints proves to be challenging, as they are often nonconvex and hard to evaluate without resolving to computationally expensive sampling techniques. These difficulties have motivated further research into tractable forms of chance constraints, which led to several convex approximations [21]. This work will focus on using Boole's inequality for constraint relaxation. A detailed description in the context of trajectories will follow in Section II-B.

A. Problem Formulation

Consider the constrained optimal control problem with probabilistic state and action constraints and unknown stochastic time-discrete transition dynamics

$$\begin{array}{ll} \max_{\boldsymbol{A}} & J(\boldsymbol{s}, \boldsymbol{A}), \\ \text{s.t.} & \boldsymbol{s}_{t+1} \sim \mathcal{P}(\boldsymbol{s}_{t+1} | \boldsymbol{s}_t, \boldsymbol{a}_t), \\ & \Pr(\boldsymbol{s}_{0:T} \in \mathcal{S}) \geq 1 - \theta, \\ & \Pr(\boldsymbol{a}_{0:T-1} \in \mathcal{A}) \geq 1 - \vartheta, \end{array}$$

where S and A are the feasible state and action spaces respectively. The probability levels θ , ϑ are hyperparameters that influence the risk behavior in terms of violating the constraints. The goal of this constrained optimization is to maximize the objective by finding the optimal action sequence A. In general, we consider the expected cumulative reward for a trajectory of length T in the quadratic form

$$J(\boldsymbol{s}, \boldsymbol{A}) = -\mathbb{E}\Big[\sum_{t=0}^{T-1} (\boldsymbol{s}_t - \boldsymbol{s}_{g,t})^{\mathsf{T}} \boldsymbol{M}_t (\boldsymbol{s}_t - \boldsymbol{s}_{g,t}) + \boldsymbol{a}_t^{\mathsf{T}} \boldsymbol{D}_t \boldsymbol{a}_t + (\boldsymbol{s}_T - \boldsymbol{s}_{g,T})^{\mathsf{T}} \boldsymbol{M}_T (\boldsymbol{s}_T - \boldsymbol{s}_{g,T})\Big], \qquad (1)$$

where M and D are positive-definite weight matrices of appropriate dimensions and s_g is the target state. Note that a quadratic objective is not necessarily required, and non-quadratic objectives can be locally approximated by quadratic forms.

B. Relaxation of Chance Constraints

Chance constraints can be conservatively relaxed by applying Boole's inequality [22]–[24]. For the purpose of brevity, only upper-bound state constraints are considered. However, the same relaxation procedure can be straightforwardly applied to obtain a lower-bound and to relax the action constraints. Generally, the state-linear joint chance constraint for a whole trajectory is formulated as

$$\Pr(\boldsymbol{s}_{0:T} \in \mathcal{S}) = \Pr(\bigcap_{t=0}^{T} \boldsymbol{s}_{t} \in \mathcal{S}) \ge 1 - \theta,$$
$$= \Pr(\bigcap_{t=0}^{T} \boldsymbol{h}_{t}^{\mathsf{T}} \boldsymbol{s}_{t} \le b_{t}) \ge 1 - \theta.$$
(2)

where h_t and b_t parameterize the half-plane defined by the constraints. Consequently, the probability of a trajectory to be within a feasible set is constrained to be higher than a probability $1 - \theta$. In the framework of stochastic programming, it is usually beneficial to reformulate Equation (2) into separate inequalities over individual constraints [20], which is achieved by transforming the intersection operator into a union operator according to rules of probability.

$$\Pr(\bigcap_{t=0}^{T} \boldsymbol{h}_{t}^{\mathsf{T}} \boldsymbol{s}_{t} \leq b_{t}) = 1 - \Pr(\bigcup_{t=0}^{T} \boldsymbol{h}_{t}^{\mathsf{T}} \boldsymbol{s}_{t} > b_{t}),$$
$$\geq 1 - \sum_{t=0}^{T} 1 - \Pr(\boldsymbol{h}_{t}^{\mathsf{T}} \boldsymbol{s}_{t} \leq b_{t}). \quad (3)$$

The sum in Inequality (3) can now be treated as a collection of single probabilities per time-step

$$\sum_{t=0}^{T} 1 - \Pr(\boldsymbol{h}_{t}^{\mathsf{T}} \boldsymbol{s}_{t} \leq b_{t}) \leq \boldsymbol{\theta},$$
$$\Pr(\boldsymbol{h}_{t}^{\mathsf{T}} \boldsymbol{s}_{t} \leq b_{t}) \geq 1 - \boldsymbol{\theta}_{t}, \tag{4}$$

where $\sum_{t=0}^{T} \theta_t = \theta$. By assuming a Gaussian probability density, a common assumption in control applications, Equation (4) is rewritten using the cumulative density function

$$\frac{1}{2} \left[1 + \operatorname{erf}\left(\frac{b_t - \boldsymbol{h}_t^{\mathsf{T}} \boldsymbol{\mu}_{\boldsymbol{s}_t}}{\sqrt{2\boldsymbol{h}_t^{\mathsf{T}} \boldsymbol{\Sigma}_{\boldsymbol{s}_t} \boldsymbol{h}_t}}\right) \right] \ge 1 - \theta_t,$$

$$b_t - \boldsymbol{h}_t^{\mathsf{T}} \boldsymbol{\mu}_{\boldsymbol{s}_t} - \sqrt{2\boldsymbol{h}_t^{\mathsf{T}} \boldsymbol{\Sigma}_{\boldsymbol{s}_t} \boldsymbol{h}_t} \operatorname{erf}^{-1}(1 - 2\theta_t) \ge 0, \quad (5)$$

where μ_{st} and Σ_{st} are the state mean and covariance respectively. Moreover, due to properties of the error function, the inequality $\sum_{t=0}^{T} \theta_t \leq \theta < 0.5$ is conservatively enforced by setting $\theta_t = \theta/T$ and requiring $\theta < 0.5$, as in [24].

C. Iterative Linear Quadratic Gaussian Control (iLQG)

We base our trajectory optimization technique on DDP/iLQG methods. This section provides a short overview on the principles of DDP [8] and iLQG [2]. For any arbitrary time-index reward function R_t , the trajectory optimization objective is the expected cumulative reward

$$J(\boldsymbol{s}, \boldsymbol{A}) = \mathbb{E}\left[\sum_{t=0}^{T-1} R_t(\boldsymbol{s}_t, \boldsymbol{a}_t) + R_T(\boldsymbol{s}_T)\right]$$

DDP and iLQG leverage the principle of dynamic programming to simplify the optimization over a complete sequence of actions $a_{0:T-1}$ to an optimization over single actions a_t for each time-step. For this purpose the time-indexed statevalue function is introduced

$$V_t(s) = \max_{a_t} \left[R_t(s_t, a_t) + \sum_{s_{t+1}} V_{t+1}(s_{t+1}) \mathcal{P}(s_{t+1}|s_t, a_t) \right]$$

over which the dynamic programming backward recursion is performed. By assuming linear transitions dynamics and a quadratic rewards along a nominal trajectory, optimal feedback controllers can be derived in closedform. DDP and iLQG consider the perturbed state-actionvalue function $Q_t(\delta s, \delta a) = R_t(s_t + \delta s, a_t + \delta a) - R_t(s_t, a_t) + V_{t+1} (\mathcal{P}(s_t + \delta s, a_t + \delta a)) - V_{t+1} (\mathcal{P}(s_t, a_t))$, resulting from a second order Taylor approximation

$$Q_t(\delta \boldsymbol{s}, \delta \boldsymbol{a}) \approx \frac{1}{2} \begin{bmatrix} 1\\ \delta \boldsymbol{s}\\ \delta \boldsymbol{a} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} 0 & \boldsymbol{Q}_{\boldsymbol{s},t}^{\mathsf{T}} & \boldsymbol{Q}_{\boldsymbol{a},t}^{\mathsf{T}}\\ \boldsymbol{Q}_{\boldsymbol{s},t} & \boldsymbol{Q}_{\boldsymbol{s}\boldsymbol{s},t} & \boldsymbol{Q}_{\boldsymbol{s}\boldsymbol{a},t} \\ \boldsymbol{Q}_{\boldsymbol{a},t} & \boldsymbol{Q}_{\boldsymbol{a}\boldsymbol{s},t} & \boldsymbol{Q}_{\boldsymbol{a}\boldsymbol{a},t} \end{bmatrix} \begin{bmatrix} 1\\ \delta \boldsymbol{s}\\ \delta \boldsymbol{a} \end{bmatrix}.$$

The subscripts s and a stand for the first and second order approximations. The entries of $Q_t(\delta s, \delta a)$ are given as

$$\begin{split} \boldsymbol{Q}_{s,t} &= \boldsymbol{R}_{s,t} + \mathcal{P}_{s,t}^{\mathsf{T}} \boldsymbol{V}_{s,t+1}, \\ \boldsymbol{Q}_{a,t} &= \boldsymbol{R}_{a,t} + \mathcal{P}_{a,t}^{\mathsf{T}} \boldsymbol{V}_{s,t+1}, \\ \boldsymbol{Q}_{ss,t} &= \boldsymbol{R}_{ss,t} + \mathcal{P}_{s,t}^{\mathsf{T}} \boldsymbol{V}_{ss,t+1} \mathcal{P}_{s,t} + \boldsymbol{V}_{s,t+1} \mathcal{P}_{ss,t}, \\ \boldsymbol{Q}_{aa,t} &= \boldsymbol{R}_{aa,t} + \mathcal{P}_{a,t}^{\mathsf{T}} \boldsymbol{V}_{ss,t+1} \mathcal{P}_{a,t} + \boldsymbol{V}_{s,t+1} \mathcal{P}_{aa,t}, \\ \boldsymbol{Q}_{as,t} &= \boldsymbol{R}_{as,t} + \mathcal{P}_{a,t}^{\mathsf{T}} \boldsymbol{V}_{ss,t+1} \mathcal{P}_{s,t} + \boldsymbol{V}_{s,t+1} \mathcal{P}_{as,t}. \end{split}$$

The main difference of iLQG compared to DDP is in neglecting the second order derivatives of the dynamics in iLQG. Given these approximations the optimal feedback controller is given as $\delta a^* = -Q_{aa,t}^{-1}(Q_a + Q_{as,t}\delta s) = K_t \delta s + k_t$. Inserting δa^* into $Q_t(\delta s, \delta a)$ returns the update equations of the state-value function per time-step

$$\Delta V_t = -rac{1}{2} Q_{a,t} Q_{aa,t}^{-1} Q_{a,t}, \ V_{s,t} = Q_{s,t} - Q_{a,t} Q_{aa,t}^{-1} Q_{as,t}, \ V_{ss,t} = Q_{ss,t} - Q_{sa,t} Q_{aa,t}^{-1} Q_{as,t},$$

During the forward pass, new trajectories of the stochastic non-linear dynamics are sampled by propagating the actions through the real system

$$a_{t} = a_{r,t} + k_{t} + K_{t}(s_{t} - s_{r,t}),$$

$$s_{t+1} \sim \mathcal{P}(s_{t+1}|s_{t}, a_{t}), \quad s_{0} = s_{r,0},$$
(6)

where $s_{r,t}$, $a_{r,t}$ denote the mean state and action at time t from the last iteration and are also referred to as the nominal or reference trajectory, here denoted by the subscript r.

Special care has to be taken during the backward pass of DDP and iLQG to ensure that $Q_{aa,t}$ is negative-definite, which has inspired different regularization schemes. In DDP, this regularization is commonly applied to $Q_{aa,t}$ as $\tilde{Q}_{aa,t} = Q_{aa,t} - \mu I$, with $\mu \ge 0$. However, other regularizations directly affecting the value function have been shown to be more effective [2], and will be used throughout this work.

D. Augmented Linearized Closed-Loop System

To formulate the chance-constrained optimization problem, we first introduce the notation and system description of the online-fitted time-variant linear system. Following [19], our approach optimizes the feedforward terms of the control, while satisfying the constraints for the linearized dynamics and maintains the feedback gains computed during the backward pass of DDP/iLQG.

Given N trajectories from the non-linear system as described in Equation (6), we fit linear-Gaussian models to the sampled data via regularized linear regression. Consequently we obtain the transition and control matrices A_t, B_t , as well as the bias vector c_t for each time-step. The resulting timevariant linear dynamics $s_{t+1} = A_t s_t + B_t a_t + c_t + w_t$, with $w_t \sim \mathcal{N}(\mathbf{0}, \Sigma_t)$, and the controller $a_t = K_t(s_t - s_{r,t}) + k_t + a_{r,t}$ are used to formulate the closed-loop linear system $s_{t+1} = \hat{A}_t s_t + B_t k_t + d_t + w_t$, where $\hat{A}_t = A_t + B_t K_t$ and $d_t = c_t - B_t K_t s_{r,t} + B_t a_{r,t}$.

To represent the closed-loop system over an entire trajectory we use the augmented notation

$$egin{aligned} & ilde{s} = egin{bmatrix} s_0 \ s_1 \ dots \ s_T \end{bmatrix}, & ilde{k} = egin{bmatrix} k_0 \ k_1 \ dots \ k_T \ dots \ k_T \ dots \ k_T \ dots \ \ \ \ \ \ \ \$$

The augmented weighting matrices for the quadratic objective take the form

$$\begin{split} \bar{\boldsymbol{M}} &= \operatorname{diag}(\boldsymbol{M}_0, \dots, \boldsymbol{M}_T), \ \bar{\boldsymbol{D}} &= \operatorname{diag}(\boldsymbol{D}_0, \dots, \boldsymbol{D}_{T-1}), \\ \tilde{\boldsymbol{M}}_C &= \operatorname{diag}(\boldsymbol{M}_0 + \boldsymbol{K}_0^{\mathsf{T}} \boldsymbol{D}_0 \boldsymbol{K}_0, \dots, \\ \boldsymbol{M}_{T-1} + \boldsymbol{K}_{T-1}^{\mathsf{T}} \boldsymbol{D}_{T-1} \boldsymbol{K}_{T-1}, \boldsymbol{M}_T), \\ \tilde{\boldsymbol{K}} &= \operatorname{diag}(\boldsymbol{K}_0, \dots, \boldsymbol{K}_{T-1}), \end{split}$$

and the closed-loop linearized stochastic dynamics is written in terms of the augmented notation as

$$\tilde{s} = \tilde{A}s_0 + \tilde{B}\tilde{k} + \tilde{G}\tilde{w} + \tilde{G}\tilde{d}, \qquad (7)$$

which in turn can be decomposed to the mean and covariance of a Gaussian state density

$$egin{aligned} egin{aligned} \mu_{ ilde{s}} &= ilde{A} m{s}_0 + ilde{B} m{k} + ilde{G} m{d}, \ m{ ilde{\Sigma}}_{ ilde{s}} &= ilde{A} \Sigma_{m{s}_0} m{ ilde{A}}^{ op} + m{ ilde{G}} m{ ilde{\Sigma}}_{m{ ilde{w}}} m{ ilde{G}}^{ op} \end{aligned}$$

where $\tilde{\Sigma}_{\tilde{w}}$ are the stacked estimates of the covariance for each time-step, taken under the N samples drawn during the

last forward pass. Furthermore, given the feedback gains, we compute the action covariance along the trajectory

$$ilde{\Sigma}_{ ilde{m{a}}} = ilde{m{K}} ilde{m{A}} \Sigma_{m{s}_0} ilde{m{A}}^{^{\intercal}} ilde{m{K}}^{^{\intercal}} + ilde{m{K}} ilde{m{G}} ilde{m{\Sigma}}_{ ilde{m{w}}} ilde{m{G}}^{^{\intercal}} ilde{m{K}}^{^{\intercal}}$$

E. Augmented Objective and Relaxed Chance Constraints

We simplify Objective (1) by using the stacked notation and the closed-loop matrices from Section II-D

$$\begin{split} J(\tilde{s}, \tilde{a}) &= -\mathbb{E}[\tilde{s}^{\mathsf{T}} \tilde{M}_{C} \tilde{s}] + \mathbb{E}[2 \tilde{s}_{g}^{\mathsf{T}} \tilde{M} \tilde{s}] - \mathbb{E}[\tilde{s}_{g}^{\mathsf{T}} \tilde{M} \tilde{s}_{g}] \dots \\ \dots &+ \mathbb{E}[2 \tilde{s}_{r}^{\mathsf{T}} \tilde{K}^{\mathsf{T}} \tilde{D} \tilde{K} \tilde{s}] - \mathbb{E}[2 \tilde{a}_{r}^{\mathsf{T}} \tilde{D} \tilde{K} \tilde{s}] - \mathbb{E}[2 \tilde{k}^{\mathsf{T}} \tilde{D} \tilde{K} \tilde{s}] \dots \\ \dots &- \mathbb{E}[\tilde{s}_{r}^{\mathsf{T}} \tilde{K}^{\mathsf{T}} \tilde{D} \tilde{K} \tilde{s}_{r}] + \mathbb{E}[2 \tilde{a}_{r}^{\mathsf{T}} \tilde{D} \tilde{K} \tilde{s}_{r}] + \mathbb{E}[2 \tilde{k}^{\mathsf{T}} \tilde{D} \tilde{K} \tilde{s}_{r}] \dots \\ \dots &- \mathbb{E}[\tilde{a}_{r}^{\mathsf{T}} \tilde{D} \tilde{a}_{r}] - \mathbb{E}[2 \tilde{k}^{\mathsf{T}} \tilde{D} \tilde{a}_{r}] - \mathbb{E}[\tilde{k}^{\mathsf{T}} \tilde{D} \tilde{k}]. \end{split}$$

Given that the expectations are of linear-quadratic quantities under Gaussian densities, it is possible to evaluate this objective in closed-form. This objective depends only on the forward terms \tilde{k} and can be reformulated as $\tilde{J}(\tilde{k})$.

Following the relaxation presented in Section II-B and using the stacked notation we can write the upper and lower state-linear chance constraints as

$$\tilde{\boldsymbol{b}}_{u} - \tilde{\boldsymbol{h}}_{u}^{\mathsf{T}} \boldsymbol{\mu}_{\tilde{\boldsymbol{s}}} - \sqrt{2 \tilde{\boldsymbol{h}}_{u}^{\mathsf{T}} \tilde{\boldsymbol{\Sigma}}_{\tilde{\boldsymbol{s}}} \tilde{\boldsymbol{h}}_{u}} \odot \operatorname{erf}^{-1}(1 - 2 \tilde{\boldsymbol{\theta}}_{u}) \ge \boldsymbol{0}, \quad (8)$$

$$-\tilde{\boldsymbol{b}}_{l}+\tilde{\boldsymbol{h}}_{l}^{\mathsf{T}}\boldsymbol{\mu}_{\tilde{\boldsymbol{s}}}+\sqrt{2\tilde{\boldsymbol{h}}_{l}^{\mathsf{T}}\tilde{\boldsymbol{\Sigma}}_{\tilde{\boldsymbol{s}}}\tilde{\boldsymbol{h}}_{l}}\odot\operatorname{erf}^{-1}(2\tilde{\boldsymbol{\theta}}_{l}-1)\geq\mathbf{0},\quad(9)$$

where \tilde{h} and \tilde{b} parameterize the upper and lower halfplanes of the state constraints and $\tilde{\theta}_u$ and $\tilde{\theta}_l$ denote the probability values per time-step, all stacked and indexed by u and l respectively. Analogously, the action constraints of the closed-loop system can be formulated

$$\tilde{\boldsymbol{z}}_{u} - \tilde{\boldsymbol{f}}_{u}^{\mathsf{T}}(\tilde{\boldsymbol{K}}(\boldsymbol{\mu}_{\tilde{\boldsymbol{s}}} - \tilde{\boldsymbol{s}}_{r}) + \tilde{\boldsymbol{a}}_{r} + \tilde{\boldsymbol{k}}) - \boldsymbol{\lambda}_{u} \ge \boldsymbol{0}, \qquad (10)$$

$$-\tilde{z}_l + \tilde{f}_l^{\mathsf{T}}(\tilde{K}(\mu_{\tilde{s}} - \tilde{s}_r) + \tilde{a}_r + \tilde{k}) + \lambda_l \ge 0, \quad (11)$$

where $\lambda_u = \sqrt{2\tilde{f}_u^{\mathsf{T}}\tilde{\Sigma}_{\tilde{a}}\tilde{f}_u} \odot \operatorname{erf}^{-1}(1 - 2\tilde{\vartheta}_u)$ and $\lambda_l = \sqrt{2\tilde{f}_l^{\mathsf{T}}\tilde{\Sigma}_{\tilde{a}}\tilde{f}_l} \odot \operatorname{erf}^{-1}(2\tilde{\vartheta}_l - 1)$, \tilde{f} and \tilde{z} are the stacked halfplane parameters of the action constraints and $\tilde{\vartheta}_u, \tilde{\vartheta}_l$ are the stacked upper and lower bound probabilities per time-step. The operator \odot denotes the element-wise multiplication.

F. Chance-Constrained Trajectory Optimization

Based on the formulations introduced in Section II-D and Section II-E, it is possible to construct an optimization problem around the reference trajectory to find a sequence of feedforward terms \tilde{k} that maintain the Constraints (8-11).

The resulting optimization is a quadratic program with linear constraints in \tilde{k} . Thus, the probabilistic problem reduces to a deterministic one, which can be solved efficiently with many numerical solvers, for example, qpOASES [25] within the CasADi framework [26]. The complete dynamic programming and optimization loop is described in Algorithm 1 and is summarized as follows: During an initial forward pass, we obtain N trajectory samples, around which the dynamics is linearized for each time-step. The linearized dynamics is used to perform the backward pass of iLQG and obtain the feedback and feedforward controllers along

Algorithm 1 Chance-Constrained Trajectory Opt. (CCTO)

Input: $\theta_{u,t}, \theta_{l,t}, \vartheta_{u,t}, \vartheta_{l,t}, \alpha, N$ Output: $K_t, k_t, s_{r,t}, a_{r,t}$ 1: $a_t^{1:N}, s_t^{1:N} \leftarrow$ forwardPass $(a_{r,t}, s_{r,t}, K_t, k_t, \alpha)$ 2: while not converged do 3: $a_{r,t}, s_{r,t} \leftarrow$ meanTraj $(a_t^{1:N}, s_t^{1:N})$ 4: $A_t, B_t, c_t \leftarrow$ fitDynamics $(a_t^{1:N}, s_t^{1:N})$ 5: $K_t, k_t^* \leftarrow$ backwardPass (A_t, B_t) 6: $k_t \leftarrow$ solveQP $(A_t, B_t, c_t, K_t, k_t^*, \theta_{u,t}, \theta_{l,t}, \vartheta_{u,t}, \vartheta_{l,t})$ 7: $a_t^{1:N}, s_t^{1:N} \leftarrow$ forwardPass $(a_{r,t}, s_{r,t}, K_t, k_t, \alpha)$ 8: end while

the reference trajectory. These controllers are then used to formulate the closed-loop linearized system with the stacked notation and to warm-start the quadratic program. The solution of the constrained program returns the optimal feedforward sequence k_t , which is used to perform the next forward pass and linearization. Following [2], we also use the hyperparameter α that scales the feedforward control in order to keep the next forward pass of the non-linear system in a valid trust-region around the linear-quadratic approximations.

III. EMPIRICAL EVALUATION

We evaluate our approach on two highly non-linear dynamical tasks, the Furuta pendulum [27] and a Cart-Pole environment. Both problems are under-actuated and have state and actions constraints. We consider quadratic reward functions for both experiments and set the probability values for violating the constraints to $\theta_u = \theta_l = \vartheta_u = \vartheta_l = 0.01$.

a) Furuta Pendulum Swing-Up: In the Furuta pendulum the state is represented by the angles of both links and the corresponding angular velocities. Only the horizontal link is actuated and is subject to both state and the action constraints. To make the environment stochastic, we introduce both action and process noise. We run our experiment under identical conditions for CCTO and iLOG. We fix the feedforward scalar α to 0.05 for both algorithms and perform 20 seeded trials, each with 45 iterations, 50 rollouts per iteration. The resulting performance curve of both algorithms can be seen in Figure 1. Furthermore, we present the planned nominal trajectories, as well as the planned nominal actions of both algorithms for one trial. The filled space is the area between the minimum and maximum values of states and actions and should not be confused with a probability distribution over trajectories. The advantage of our approach is clear. CCTO reaches better overall performance with a higher final reward and smaller standard deviation, Table I. iLQG plans frequently and consistently to violate the constraints, while CCTO keeps the state and action trajectories within a feasible space. This consideration leads to an improved approximation of the nonlinear system dynamic and allows CCTO to perform robust improvement steps during the optimization. This result is affirmed by the low regularization values of CCTO, Table II.

b) Cart-Pole Swing-Up: For the well-known Cart-Pole environment, we consider constraints on the position of the cart as well as on the action. To make the task more



Fig. 1: Total-reward curve reflecting the performance of iLQG and CCTO for the Furuta pendulum swing-up task (left). In addition, we show the space (min. and max.) of planned nominal trajectories of the constrained angle (middle) and the corresponding executed actions (right), CCTO (blue), iLQG (red). CCTO obeys the physical limits of the system, while iLQG drives the dynamics against the constraints (green). These violations lead to poor linear approximations of the dynamics and an overall slightly lower mean and higher variance performance of iLQG.

Iteration	10	30	45
ССТО	$-6.8(\pm 0.32)$	$-1.3(\pm 0.11)$	$-0.65(\pm 0.6)$
iLQG	$-4.3(\pm 0.46)$	$-1.6(\pm 0.39)$	$-1.1(\pm 0.53)$

TABLE I: Mean total reward and standard deviation of the Furuta swing-up task scaled by 1e-2.

Iteration	10	30	45
ССТО	0	$\mathbf{2.5e}{-8}$	1e-4
iLQG	0	2.85e38	5e80

TABLE II: Mean regularization in the Furuta task over all trials for different iterations. CCTO needs less regularization due to avoidance of hard non-linearities.

challenging, we again apply action and process noise, enforce harsh action constraints and limit the time horizon to 100 time steps, the equivalent of 2 seconds. We evaluate iLQG and CCTO on 20 seeded trials, each with 55 iterations and 50 rollouts per iteration. We set the feedforward scaling parameter α to 0.1. Analogously to the last experiment, Figure 2 depicts the performance curve of iLQG and CCTO, as well as the spaces of planned nominal trajectories for the cart's position and the corresponding actions. In this experiment, iLQG moves very quickly towards a local optimum and does not manage to swing the Cart-Pole up. In contrast, CCTO performs the swing-up by finding a suitable nominal trajectory in the feasible constrained space. Tables III and IV reflect the performance discrepancy between both algorithms, in terms of total rewards and needed regularization.

IV. CONCLUSION AND FUTURE RESEARCH

We have proposed a new trajectory optimization technique, based on the framework of differential dynamic program-

Iteration	20	30	55
CCTO	$-2.3(\pm 0.32)$	$-1.2(\pm 0.32)$	$-0.31(\pm 0.06)$
iLQG	$-9.3(\pm 0.10)$	$-9.3(\pm 0.10)$	$-9.3(\pm 0.10)$

TABLE III: Mean total reward and standard deviation of the Cart-Pole swing-up task scaled by 1e-2.

Iteration	20	30	55
CCTO	0	0	0
iLQG	5.7e39	1e80	1e80

TABLE IV: Mean regularization in the Cart-Pole task over all trials for different iterations. CCTO needs less regularization due to avoidance of hard non-linearities.

ming, that takes into consideration probabilistic chance constraints in stochastic environments with unknown non-linear dynamics. We used Boole's inequality to conservatively relax the non-convex chance constraints, enabling us to formulate a constrained quadratic program and optimize the nominal trajectory to stay inside the feasible set defined by the probabilistic linear state and action limits. We have provided a thorough derivation of our approach and empirically demonstrated the advantage of enforcing physical limits on two simulated highly dynamical and stochastic non-linear systems. The results indicate that incorporating the chance constraints leads to higher fidelity in the online-fitted local linear-quadratic approximations, and consequently greatly influences the robustness of the iterative optimization process. This observation is reflected in very low regularizations in comparison to standard iLQG.

In future research, we will extend our optimization to include not only the nominal trajectory but also the feedback gains, and we will consider optimizing the probabilistic



Fig. 2: Total-reward curve reflecting the performance of iLQG and CCTO for the Cart-Pole swing-up task (left). Furthermore, we show the space (min. and max.) of planned nominal trajectories of the constrained position (middle) and the corresponding executed actions (right), CCTO (blue), iLQG (red). CCTO obeys the physical limits of the system, while iLQG drives the dynamics against the constraints (green). These violations, especially those of the action constraint cause iLQG to get stuck in a poor local optimum, while CCTO is able to solve the task and perform the swing-up.

constraint bounds via risk allocation to achieve dynamic risk measures across time and iterations. In addition, we plan to move to the fully stochastic optimization framework of maximum-entropy iLQG [6] to avoid regularization heuristics of the DDP framework.

REFERENCES

- [1] M. P. Deisenroth, G. Neumann, and J. Peters, "A survey on policy search for robotics," *Foundations and Trends*(R) *in Robotics*, 2013.
- [2] Y. Tassa, T. Erez, and E. Todorov, "Synthesis and stabilization of complex behaviors through online trajectory optimization," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012.
- [3] M. P. Deisenroth, D. Fox, and C. E. Rasmussen, "Gaussian processes for data-efficient learning in robotics and control," *IEEE transactions* on pattern analysis and machine intelligence, 2015.
- [4] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *The Journal of Machine Learning Research*, 2016.
- [5] M. Deisenroth and C. E. Rasmussen, "PILCO: A model-based and data-efficient approach to policy search," in *Proceedings of the 28th International Conference on machine learning*, 2011.
- [6] S. Levine and P. Abbeel, "Learning neural network policies with guided policy search under unknown dynamics," in Advances in Neural Information Processing Systems, 2014.
- [7] H. Abdulsamad, O. Arenz, J. Peters, and G. Neumann, "Stateregularized policy search for linearized dynamical systems," in *International Conference on Automated Planning and Scheduling*, 2017.
- [8] D. H. Jacobson and D. Q. Mayne, "Differential dynamic programming," 1970.
- [9] E. Todorov and W. Li, "A generalized iterative LQG method for locally-optimal feedback control of constrained nonlinear stochastic systems," in *American Control Conference*. IEEE, 2005.
- [10] B. D. Anderson and J. B. Moore, *Optimal control: Linear quadratic methods*. Courier Corporation, 2007.
- [11] O. Von Stryk and R. Bulirsch, "Direct and indirect methods for trajectory optimization," Annals of operations research, 1992.
- [12] Y. Tassa, N. Mansard, and E. Todorov, "Control-limited differential dynamic programming," in *IEEE International Conference on Robotics* and Automation. IEEE, 2014.

- [13] B. Plancher, Z. Manchester, and S. Kuindersma, "Constrained unscented dynamic programming," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2017.
- [14] Z. Xie, C. K. Liu, and K. Hauser, "Differential dynamic programming with nonlinear constraints," in *IEEE International Conference on Robotics and Automation*. IEEE, 2017.
- [15] J. Van Den Berg, S. Patil, and R. Alterovitz, "Efficient approximate value iteration for continuous gaussian pomdps," in AAAI Conference on Artificial Intelligence, 2012.
- [16] D. H. Van Hessem, "Stochastic inequality constrained closed-loop model predictive control–with application to chemical process operation," 2004.
- [17] L. Blackmore and M. Ono, "Convex chance constrained predictive control without sampling," in AIAA Guidance, Navigation, and Control Conference, 2009.
- [18] M. P. Vitus and C. J. Tomlin, "Closed-loop belief space planning for linear, gaussian systems," in *IEEE International Conference on Robotics and Automation*. IEEE, 2011.
- [19] G. Kurz, M. Dolgov, and U. D. Hanebeck, "Progressive closed-loop chance-constrained control," in *International Conference on Information Fusion*. IEEE, 2016.
- [20] A. Prékopa, Stochastic programming. Springer Science & Business Media, 2013.
- [21] A. Nemirovski and A. Shapiro, "Convex approximations of chance constrained programs," SIAM Journal on Optimization, 2006.
- [22] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [23] M. P. Vitus and C. J. Tomlin, "On feedback design and risk allocation in chance constrained control," in *IEEE Conference on Decision and Control and European Control Conference*. IEEE, 2011.
- [24] S. Jha and V. Raman, "On optimal control of stochastic linear hybrid systems," in *International Conference on Formal Modeling and Analysis of Timed Systems*. Springer, 2016.
- [25] H. J. Ferreau, C. Kirches, A. Potschka, H. G. Bock, and M. Diehl, "qpOASES: A parametric active-set algorithm for quadratic programming," *Mathematical Programming Computation*, 2014.
- [26] J. A. Andersson, J. Gillis, G. Horn, J. B. Rawlings, and M. Diehl, "Casadi: a software framework for nonlinear optimization and optimal control," *Mathematical Programming Computation*, 2018.
- [27] K. Furuta, M. Yamakita, and S. Kobayashi, "Swing-up control of inverted pendulum using pseudo-state feedback," *Journal of Systems* and Control Engineering, 1992.