



# Bayesian optimization with safety constraints: safe and automatic parameter tuning in robotics

Felix Berkenkamp<sup>1</sup> · Andreas Krause<sup>1</sup> · Angela P. Schoellig<sup>2</sup>

Received: 2 May 2020 / Revised: 2 May 2020 / Accepted: 3 June 2021  
© The Author(s) 2021

## Abstract

Selecting the right tuning parameters for algorithms is a prevalent problem in machine learning that can significantly affect the performance of algorithms. Data-efficient optimization algorithms, such as Bayesian optimization, have been used to automate this process. During experiments on real-world systems such as robotic platforms these methods can evaluate unsafe parameters that lead to safety-critical system failures and can destroy the system. Recently, a safe Bayesian optimization algorithm, called SAFE<sub>OPT</sub>, has been developed, which guarantees that the performance of the system never falls below a critical value; that is, safety is defined based on the performance function. However, coupling performance and safety is often not desirable in practice, since they are often opposing objectives. In this paper, we present a generalized algorithm that allows for multiple safety constraints separate from the objective. Given an initial set of safe parameters, the algorithm maximizes performance but only evaluates parameters that satisfy safety for all constraints with high probability. To this end, it carefully explores the parameter space by exploiting regularity assumptions in terms of a Gaussian process prior. Moreover, we show how context variables can be used to safely transfer knowledge to new situations and tasks. We provide a theoretical analysis and demonstrate that the proposed algorithm enables fast, automatic, and safe optimization of tuning parameters in experiments on a quadrotor vehicle.

**Keywords** Bayesian optimization · Safety constraints · Safe exploration · Reinforcement learning · Robotics

---

Editors: Daniel Fremont, Mykel Kochenderfer, Alessio Lomuscio, Dragos Margineantu, Cheng Soon-Ong.

---

✉ Felix Berkenkamp  
fberkenkamp@gmail.com

<sup>1</sup> Department of Computer Science, ETH Zurich, Zurich, Switzerland

<sup>2</sup> Institute for Aerospace Studies, University of Toronto, Toronto, Canada

# 1 Introduction

Reinforcement learning algorithms (Sutton & Barto, 1998) have demonstrated their ability to find optimal policy parameters by interacting with a system. However, they face two obstacles in the real world: To solve tasks they often require large amounts of training data (Lillicrap et al., 2015), which is impractical on physical system. Secondly, they often operate in a trial and error process that fundamentally relies on trying out the different available actions and observing their consequences. This is not practical for real-world systems, where safety and the ability to operate within constraints imposed by an environment are critical prerequisites. For example, in robotics, learning algorithms often face large prior uncertainties and failures can cause serious damage to the robot and its environment (Schaal & Atkeson, 2010).

One class of algorithms that is sufficiently data-efficient for use in the real world is Bayesian optimization (Mockus, 2012). These methods learn a Bayesian model that maps tuning parameters to the resulting performance objective and use it to evaluate *informative* parameters, which leads to data-efficient optimization that can provably find optimal parameters after only few parameter evaluations (Bull, 2011; Srinivas et al., 2012). These methods have been successfully applied in experiments on real robots (Lizotte et al., 2007), but safety is not typically considered.

The SAFE<sub>OPT</sub> algorithm by Sui et al. (2015) is a Bayesian optimization algorithm that, in addition to being data-efficient, also considers safety *during* the learning process. It defines safety as a minimum performance requirement that cannot be violated when evaluating new parameters. To this end, it effectively trades off finding the optimal parameters and learning about the safety of other parameters. However, defining safety in terms of performance only is often a restrictive requirement in practice, where the two properties are often conflicting. For example, the shortest path for a robot may involve flying close to an obstacle, while safety would require keeping a distance.

*Contributions* In this paper, we present an extension of the SAFE<sub>OPT</sub> algorithm (Sui et al., 2015) that considers multiple, arbitrary safety constraints decoupled from the performance objective. This generalization retains the desirable sample-efficient properties of normal Bayesian optimization, but carefully explores the parameter space in order to maximize performance while guaranteeing safety with high probability. We extend the theory of SAFE<sub>OPT</sub> to account for these additional constraints and show that similar theoretical guarantees can be obtained for the more general setting. We then relax the assumptions used in the proofs to obtain a more practical version of the algorithm and additionally show that the safety guarantees carry over to this case.

The second main contribution next to the theoretical foundation, is an extensive experimental evaluation of the method, where we consider the problem of safely optimizing linear and nonlinear laws on a quadrotor vehicle. The experiments demonstrate that the proposed approach is able to safely optimize parameters of a control algorithms while respecting safety constraints with high probability. Moreover, we show how ideas from context-based optimization (Krause & Ong, 2011) can be used to safely transfer knowledge in order to obtain environment-dependent control laws. For example, in our experiments we optimize a control law for different flying speeds of a quadrotor vehicle.

The experimental results in Sect. 5.3 were previously shown in Berkenkamp et al. (2016b) and were limited to a safety constraint on performance only and without additional theoretical results.

## 1.1 Related work

In control theory and robotics, safety is often guaranteed by designing controllers that are robust with respect to errors in a given dynamical model of the system (Zhou & Doyle, 1998). To avoid the conservatism that comes with designing robust controllers, these methods can be combined with machine learning in order to improve the model estimate while guaranteeing that the system remains stable (Berkenkamp & Schoellig, 2015; Berkenkamp et al., 2017) or never leaves a safe subset of the state space (Ostafew et al. 2016; Aswani et al. 2013; Akametalu et al. 2014). In the case when the system model is known, but safety constraints are not, Moldovan and Abbeel (2012), Turchetta et al. (2016) use planning methods to safely explore the state space. These methods either require or learn a dynamics model of the system, which can be challenging in practice. In contrast, we do not assume to have access to a model of the system, but aim to directly optimize the parameters of a control algorithm, without violating safety constraints.

The objective of learning optimal policies has been extensively studied in the reinforcement learning community (Sutton & Barto, 1998). In particular, the area of policy search considers the objective of optimizing the parameters of control algorithms (Kober & Peters, 2013). The state of the art methods are based on estimating the gradients of the performance function (Peters & Schaal, 2006, 2008). As a result, typically multiple evaluations of very similar parameters are conducted on the real system in order to estimate the gradients, which means the approaches are often not data-efficient. Dalal et al. (2018) combine these approaches with a known model to ensure safety at the next time step. Without a known model, safety in gradient-based policy search has been considered by disallowing exploration into areas of the state space that have not been visited before (Garcia & Fernández, 2012; Hans et al., 2008). Along those lines, Achiam et al. (2017) propose a constrained policy optimization algorithm that guarantees constraint satisfaction. In contrast, our method is gradient-free, so that it can explore the parameter space globally in a more data-efficient manner. At the same time, we provide high-probability worst-case guarantees for not violating safety constraints during the optimization process.

Bayesian optimization algorithms (Jones, 2001; Mockus, 2012) without safety constraints have been successfully applied to robotic tasks before. For example, Calandra et al. (2014a), Lizotte et al. (2007) use Bayesian optimization to optimize the gait of legged robots and Tesch et al. (2011) optimize the control parameters of a snake-like robot. Marco et al. (2017) optimize the weighting matrices of an LQR controller for an inverted pendulum by exploiting additional information from a simulator and Calandra et al. (2014b) compares several different Bayesian optimization methods for bipedal locomotion. While these examples illustrate the potential of Bayesian optimization methods in robotics, none of these examples explicitly considers safety as a requirement.

While with known constraints one could constrain the search domain of Bayesian optimization to the feasible region [as in convex optimization (Boyd et al., 2004)], for unknown constraints we need to actively learn about the *a priori* unknown constraint. Gelbart et al. (2014) introduce an algorithm to optimize an unknown function subject to an unknown constraint. However, this constraint was not considered to be safety-critical; that is, the algorithm is allowed to evaluate unsafe parameters. Schreiter et al. (2015) consider safety constraints for active learning, where the goal is to identify parameters that are safe, rather than optimizing the process directly, while Sui et al. (2015) present a similar algorithm to safely optimize an objective function.

Since the results in this paper were made available in Berkenkamp et al. (2016a), there has been significant follow-up work. Sui et al. (2018) uses our model with multiple constraints, but considers a two-stages algorithm that first learns about the safe set and then aims to find the optimal parameters within this safe set. Turchetta et al. (2019) instead avoids learning about the safe set explicitly, and only explicitly learns about the safety of actions when it is required to optimize the performance objective.

## 2 Problem statement

We consider a given algorithm that is used to accomplish a certain task with a robot. In general, this algorithm is arbitrary and may contain several components including vision, state estimation, planning, and control laws. The algorithm depends on tuning parameters  $\mathbf{a} \in \mathcal{A}$  in some specified, domain  $\mathcal{A} \subseteq \mathbb{R}^d$ .

The goal is to find the parameters within  $\mathcal{A}$  that maximize a given, scalar performance measure,  $f$ . For example, this performance measure may represent the negative tracking error of a robot (Berkenkamp et al., 2016b), the average walking speed of a bipedal robot (Calandra et al., 2014a), or any other quantity that can be computed over a finite time horizon. We can only evaluate the performance measure for any parameter set  $\mathbf{a}$  on finite-time trajectories from experiments. The functional dependence of  $f$  on  $\mathbf{a}$  is not known *a priori*. In the following, we write the performance measure as a function of the parameters  $\mathbf{a}$ ,  $f : \mathcal{A} \rightarrow \mathbb{R}$ , even though measuring performance requires an experiment on the physical robot and typically depends on a trajectory of states, control inputs, and external signals.

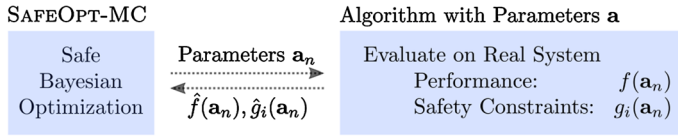
We assume that the underlying system is safety-critical; that is, there are constraints that the system must respect when evaluating parameters. Similarly to the performance measure,  $f(\mathbf{a})$ , these constraints can represent any quantity and may depend on states, inputs, or even environment variables. There are  $q$  safety constraints of the form  $g_i(\mathbf{a}) \geq 0$ ,  $g_i : \mathcal{A} \rightarrow \mathbb{R}$ ,  $i = 1 \dots q$ , which together define the safety conditions. This is without loss of generality, since any constraint function can be shifted by a constant in order to obtain this form. The functions  $g_i$  are unknown *a priori* but can be estimated through (typically noisy) experiments for a given parameter set  $\mathbf{a}$ . For example, in order to encode a state constraint on an obstacle for a robot, the safety function  $g_i(\mathbf{a})$  can return the smallest distance to the obstacle along a trajectory of states when using algorithm parameters  $\mathbf{a}$ . Note that if the functions were known in advance, we could simply exclude unsafe parameters from the set  $\mathcal{A}$ .

The overall optimization problem can be written as

$$\max_{\mathbf{a} \in \mathcal{A}} f(\mathbf{a}) \quad \text{subject to} \quad g_i(\mathbf{a}) \geq 0 \quad \forall i = 1, \dots, q. \quad (1)$$

The goal is to iteratively find the global maximum of this constrained optimization problem by, at each iteration  $n$ , selecting parameters  $\mathbf{a}_n$  and evaluating (up to noise) the corresponding function values  $f(\mathbf{a}_n)$  and  $g_i(\mathbf{a}_n)$  until the optimal parameters are found. In particular, since the constraints define the safety of the underlying system, only parameters that are inside the feasible region of (1) are allowed to be evaluated; that is, only parameters that fulfill these safety requirements on the real system (Fig. 1).

Since the functions  $f$  and  $g_i$  in (1) are unknown *a priori*, it is not generally possible to solve the corresponding optimization problem without violating the constraints. The first problem is that we do not know how to select a first, safe parameter to evaluate. In the



**Fig. 1** Overview of the algorithm. At each iteration  $n$ , the algorithm selects safe and informative parameters at which the performance and the safety constraints are evaluated. Based on the noisy information gained, the algorithm updates its belief over the functions. This safe optimization process is iterated until the safely reachable, optimal parameters are found

following, we assume that an initial safe set of parameters  $S_0 \subseteq \mathcal{A}$  is known for which the constraints are fulfilled. These serve as a starting point for the exploration of the safe region in (1). In robotics, safe initial parameters with poor performance can often be obtained from a simulation or domain knowledge.

Secondly, in order to safely explore the parameter space beyond  $S_0$ , we must be able to infer whether parameters  $\mathbf{a}$  that we have not evaluated yet are safe to use on the real system. To this end, we make regularity assumptions about the functions  $f$  and  $g_i$  in (1). We discuss these assumptions in more detail in Sect. 4.2. However, broadly speaking we make assumptions that allow us to model the functions  $f$  and  $g_i$  as a GP, construct reliable confidence intervals over the domain  $\mathcal{A}$ , and imply Lipschitz continuity properties. Using these properties, we are able to generalize safety beyond the initial, safe parameters  $S_0$ . Given the model assumptions, we require that the safety constraints hold with high probability over the entire sequence of experiments.

As a consequence of the safety requirements, it is not generally possible to find the global optimum of (1). Instead we aim to find the optimum in the part of the feasible region that is safely reachable from  $S_0$ . We formalize this precisely in Sect. 4.

Lastly, whenever we evaluate parameters on the real system, we only obtain noisy estimates of both the performance function and the constraints, since both depend on noisy sensor data along trajectories. That is, for each parameter  $\mathbf{a}$  that we evaluate, we obtain measurements  $\hat{f}(\mathbf{a}) = f(\mathbf{a}) + \omega_0$  and  $\hat{g}_i(\mathbf{a}) = g_i(\mathbf{a}) + \omega_i$ , where  $\omega_i$ ,  $i = 0, \dots, q$ , is zero-mean,  $\sigma$ -sub-Gaussian noise<sup>1</sup> (sub-Gaussian noise includes Gaussian and Uniform distributions). Note that while  $\hat{f}(\mathbf{a})$  is a random variable, we use  $\hat{f}(\mathbf{a}_n)$  to denote the measurement obtained at iteration  $n$ . In general, the noise variables may be correlated, but we do not consider this case in our theoretical analysis in Sect. 4.2. We only want to evaluate parameters where all safety constraints are fulfilled, so that  $g_i(\mathbf{a}_n) \geq 0$  for all  $i \in \{1, \dots, q\}$  and  $n \geq 1$ .

### 3 Background

In this section, we review Gaussian processes (GPs) and Bayesian optimization, which form the foundation of our safe Bayesian optimization algorithm in Sect. 4. The introduction to GPs is standard and based on Berkenkamp et al. (2016b) and Rasmussen and Williams (2006).

<sup>1</sup> A random variable  $X$  is  $\sigma$ -sub-Gaussian if its tail probability decreases exponentially fast: there exists a constant  $c$  so that  $P(|X| > t) \leq \exp(1 - ct^2/\sigma^2)$  holds for all  $t \geq 0$ . The family includes Gaussian and uniform distributions, c.f. (Eldar & Kutyniok, 2012, Ch. 5).

### 3.1 Gaussian process (GP)

Both the function  $f(\mathbf{a})$  and the safety constraints  $g_i(\mathbf{a})$  in Sect. 2 are unknown *a priori*. We use GPs as a nonparametric model to approximate these unknown functions over their common domain  $\mathcal{A}$ . In the following, we focus on a single function, the performance function. We extend this model to multiple functions in order to represent both performance and constraints in Sect. 3.1.1.

GPs are a popular choice for nonparametric regression in machine learning, where the goal is to find an approximation of a nonlinear map,  $f(\mathbf{a}) : \mathcal{A} \rightarrow \mathbb{R}$ , from an input vector  $\mathbf{a} \in \mathcal{A}$  to the function value  $f(\mathbf{a})$ . This is accomplished by assuming that the function values  $f(\mathbf{a})$ , associated with different values of  $\mathbf{a}$ , are random variables and that any finite number of these random variables have a joint Gaussian distribution (Rasmussen & Williams, 2006).

A GP is parameterized by a prior mean function and a covariance function  $k(\mathbf{a}, \mathbf{a}')$ , which defines the covariance of any two function values  $f(\mathbf{a})$  and  $f(\mathbf{a}')$ ,  $\mathbf{a}, \mathbf{a}' \in \mathcal{A}$ . The latter is also known as the kernel. In this work, the mean is assumed to be zero, without loss of generality. The choice of kernel function is problem-dependent and encodes assumptions about smoothness and rate of change of the unknown function. A review of potential kernels can be found in Rasmussen and Williams (2006) and more information about the kernels used in this paper is given in Sect. 5.

The GP framework can be used to predict the function value  $f(\mathbf{a}^*)$  for an arbitrary parameter  $\mathbf{a}^* \in \mathcal{A}$  based on a set of  $n$  past observations,  $\{\hat{f}(\mathbf{a}_i)\}_{i=1}^n$ , at the chosen parameters  $\mathcal{D}_n = \{\mathbf{a}_i\}_{i=1}^n$ . The GP model assumes that observations are noisy measurements of the true function value  $f(\mathbf{a})$ ; that is,  $\hat{f}(\mathbf{a}) = f(\mathbf{a}) + \omega$  with  $\omega \sim \mathcal{N}(0, \sigma^2)$ . Conditioned on these observations, the posterior distribution is a GP again with mean and variance

$$\mu_n(\mathbf{a}^*) = \mathbf{k}_n(\mathbf{a}^*)(\mathbf{K}_n + \mathbf{I}_n\sigma^2)^{-1}\hat{\mathbf{f}}_n, \quad (2)$$

$$\sigma_n^2(\mathbf{a}^*) = k(\mathbf{a}^*, \mathbf{a}^*) - \mathbf{k}_n(\mathbf{a}^*)(\mathbf{K}_n + \mathbf{I}_n\sigma^2)^{-1}\mathbf{k}_n^T(\mathbf{a}^*), \quad (3)$$

where  $\hat{\mathbf{f}}_n = [\hat{f}(\mathbf{a}_1), \dots, \hat{f}(\mathbf{a}_n)]^T$  is the vector of observed, noisy function values, the covariance matrix  $\mathbf{K}_n \in \mathbb{R}^{n \times n}$  has entries  $[\mathbf{K}_n]_{(i,j)} = k(\mathbf{a}_i, \mathbf{a}_j)$ ,  $i, j \in \{1, \dots, n\}$ , and the vector  $\mathbf{k}_n(\mathbf{a}^*) = [k(\mathbf{a}^*, \mathbf{a}_1), \dots, k(\mathbf{a}^*, \mathbf{a}_n)]$  contains the covariances between the new input  $\mathbf{a}^*$  and the observed data points in  $\mathcal{D}_n$ . The matrix  $\mathbf{I}_n \in \mathbb{R}^{n \times n}$  denotes the identity matrix.

#### 3.1.1 GPs with multiple outputs

So far, we have focused on GPs that model a single scalar function. In order to model not only the performance,  $f(\mathbf{a})$ , but also the safety constraints,  $g_i(\mathbf{a})$ , we have to consider multiple, possibly correlated functions. In the GP literature, these are usually treated by considering a matrix of kernel functions, which models the correlation between different functions (Álvarez et al., 2012). Here instead, we use an equivalent representation by considering a surrogate function,

$$h(\mathbf{a}, i) = \begin{cases} f(\mathbf{a}) & \text{if } i = 0 \\ g_i(\mathbf{a}) & \text{if } i \in \mathcal{I}_g, \end{cases} \quad (4)$$

which returns either the performance function or the individual safety constraints depending on the additional input  $i \in \mathcal{I}$  with  $\mathcal{I} = \{0, \dots, q\}$ , where  $\mathcal{I}_g = \{1, \dots, q\} \subset \mathcal{I}$  are the indices belonging to the constraints. The function  $h(\cdot, \cdot)$  is a single-output function and can

be modeled as a GP with scalar output over the extended parameter space  $\mathcal{A} \times \mathcal{I}$ , where we use  $\mu_n(\mathbf{a}, i)$  and  $\sigma_n(\mathbf{a}, i)$  to denote the GP's predictions over the extended space. For example, a joint kernel over the extended input domain for the performance function  $f(\mathbf{a})$  and one safety constraint  $g(\mathbf{a})$  may look like this:

$$k((\mathbf{a}, i), (\mathbf{a}', j)) = \begin{cases} \delta_{ij} k_f(\mathbf{a}, \mathbf{a}') + k_{fg}(\mathbf{a}, \mathbf{a}') & \text{if } i = 0 \\ \delta_{ij} k_g(\mathbf{a}, \mathbf{a}') + k_{fg}(\mathbf{a}, \mathbf{a}') & \text{if } i = 1, \end{cases}$$

where  $\delta_{ij}$  is the Kronecker delta. This kernel models the functions  $f(\mathbf{a})$  and  $g(\mathbf{a})$  with independent kernels  $k_f$  and  $k_g$  respectively, but also introduces a covariance function  $k_{fg}$  that additively models shared components between the two functions. By extending the training data by the extra parameter  $i$ , we can use the normal GP framework and predict function values and corresponding uncertainties using (2) and (3). When observing the function values, the index  $i$  is added to the parameter set  $\mathbf{a}$  for each observation. Including noise parameters inside the kernel allows to model noise correlation between the individual functions.

Importantly, using this surrogate function rather than the framework of Álvarez et al., (2012) enables us to lift theoretical results of Sui et al., (2015) to the more general case with multiple constraints and provide theoretical guarantees for our algorithm in Sect. 4.2.

In the setting with multiple outputs, at every iteration  $n$ , we obtain  $|\mathcal{I}| = q + 1$  measurements; one for each function. For ease of notation, we continue to write  $\mu_n$  and  $\sigma_n$ , even though we have obtained  $n \cdot (q + 1)$  measurements at locations  $\mathcal{D}_n \times \mathcal{I}$  in the extended parameter space.

### 3.2 Bayesian optimization

Bayesian optimization aims to find the global maximum of an unknown function (Mockus, 2012). The framework assumes that evaluating the function is expensive, while computational resources are relatively cheap. This fits our problem in Sect. 2, where each evaluation of the performance function corresponds to an experiment on the real system, which takes time and causes wear in the robotic system.

In general, Bayesian optimization models the objective function as a random function and uses this model to determine informative sample locations. A popular approach is to model the underlying function as a GP, see Sect. 3.1. GP-based methods use the posterior mean and variance predictions in (2) and (3) to compute the next sample location. For example, according to the GP-UCB (GP-Upper Confidence Bound) algorithm by Srinivas et al. (2012), the next sample location is

$$\mathbf{a}_n = \operatorname{argmax}_{\mathbf{a} \in \mathcal{A}} \mu_{n-1}(\mathbf{a}) + \beta_n^{1/2} \sigma_{n-1}(\mathbf{a}), \quad (5)$$

where  $\beta_n$  is an iteration-dependent scalar that reflects the confidence interval of the GP. Intuitively, (5) selects new evaluation points at locations where the upper bound of the confidence interval of the GP estimate is maximal. Repeatedly evaluating the system at locations given by (5) improves the mean estimate of the underlying function and decreases the uncertainty at candidate locations for the maximum, such that the global maximum is provably found eventually (Srinivas et al., 2012).

While (5) is also an optimization problem, its solution does not require any evaluations on the real system and only uses the GP model. This reflects the assumption of cheap computational resources. In practice, Bayesian optimization typically focuses on low-dimensional problems. However, this can be scaled up by discovering a low-dimensional

subspace of  $\mathcal{A}$  for Bayesian optimization (Djolonga et al., 2013; Wang et al., 2013) or encoding additional structure in the kernel (Duvenaud et al., 2011).

### 3.3 Contextual Bayesian optimization

Contextual Bayesian optimization is a conceptually straightforward extension of Bayesian optimization (Krause & Ong, 2011). It enables optimization of functions that depend on additional, external variables, which are called contexts. For example, the performance of a robot may depend on its battery level or the weather conditions, both of which cannot be influenced directly. Alternatively, contexts can also represent different tasks that the robot has to solve, which are specified externally by a user. The idea is to include the functional dependence on the context in the GP model, but to consider them fixed when selecting the next parameters to evaluate.

For example, given a context  $\mathbf{z} \in \mathcal{Z}$  that is fixed by the environment, we can model how the performance and constraint functions change with respect to different contexts by multiplying the kernel function  $k_a$  over the parameters, with another kernel  $k_z : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  over the contexts,

$$k((\mathbf{a}, i, \mathbf{z}), (\mathbf{a}', i', \mathbf{z}')) = k_a((\mathbf{a}, i), (\mathbf{a}', i')) \cdot k_z(\mathbf{z}, \mathbf{z}'). \quad (6)$$

This kernel structure implies that function values are correlated when both parameters and the contexts are similar. For example, we would expect selecting the same parameters  $\mathbf{a}$  for a control algorithm to lead to similar performance values if the context (e.g., the battery level) is similar.

Since contexts are not part of the optimization criterion, a modified version of (5) has to be used. It was shown by Krause and Ong (2011) that an algorithm that evaluates the GP-UCB criterion given a fixed context  $\mathbf{z}_n$ ,

$$\mathbf{a}_n = \operatorname{argmax}_{\mathbf{a} \in \mathcal{A}} \mu_{n-1}(\mathbf{a}, \mathbf{z}_n) + \beta_n^{1/2} \sigma_{n-1}(\mathbf{a}, \mathbf{z}_n), \quad (7)$$

enjoys similar convergence guarantees as normal Bayesian optimization in Sect. 3.2. Specifically, after seeing a particular context often enough, the criterion (7) will query parameters that are close-to-optimal.

### 3.4 Safe Bayesian optimization (SAFEOPT)

In this paper, we extend the safe optimization algorithm SAFEOPT (Sui et al., 2015) to multiple constraints. SAFEOPT is a Bayesian optimization algorithm, see Sect. 3.2. However, instead of optimizing the underlying performance function  $f(\mathbf{a})$  globally over the domain  $\mathcal{A}$ , it restricts itself to a safe set of parameters that achieve a certain minimum performance with high probability. This safe set is not known initially, but is estimated after each function evaluation. In this setting, the challenge is to find an appropriate evaluation strategy similar to (5), which at each iteration  $n$  not only aims to find the global maximum within the currently known safe set (exploitation), but also aims to increase the set of controllers that are known to be safe (exploration). SAFEOPT trades off between these two sets by choosing for the next experiment the parameters inside the safe set about whose performance we are the most uncertain.



## 4 SAFEPT-MC (Multiple Constraints)

In this section, we introduce the SAFEPT-MC algorithm for multiple constraints and discuss its theoretical properties. The goal of the algorithm is to solve (1) by evaluating different parameters from the domain  $\mathcal{A}$  without violating the safety constraints. To this end, any algorithm has to consider two important properties:

- (i) Expanding the region of the optimization problem that is known to be feasible or safe as much as possible without violating the constraints,
- (ii) Finding the optimal parameters within the current safe set.

For objective *i*), we need quantify the size of the safe set. To do this in a tractable manner, we focus on finite sets  $\mathcal{A}$  in the following. However, heuristic extensions to continuous domains exist (Duivenvoorden et al., 2017).

The theoretical guarantees of the algorithm rely on the continuity of the underlying function. Many commonly used kernels, such as the squared exponential (Gaussian) kernel, encode Lipschitz-continuous functions with high probability (Ghosal & Roy, 2006). We make more specific assumptions that ensure deterministic Lipschitz constants in Sect. 4.2. For now, we assume that  $f(\mathbf{a})$  and  $g_i(\mathbf{a})$  are Lipschitz continuous with Lipschitz constant  $L$  with respect to some norm<sup>2</sup>.

Since we only observe noisy estimates of both the performance function and the constraints, we cannot expect to find the entire safe region encoded by the constraints within a finite number of evaluations. Instead, we follow Sui et al. (2015) and consider learning the safety constraint up to some accuracy  $\epsilon$ . This assumption is equivalent to a minimum slack of  $\epsilon$  on the constraints in (1).

As mentioned in Sect. 2, we assume that we have access to initial, safe parameters  $S_0 \subseteq \mathcal{A}$ , for which we know that the safety constraints are satisfied *a priori*. Starting from these initial parameters, we ask what the best that any safe optimization algorithm could hope to achieve is. In particular, if we knew the safety constraint functions  $g_i(\cdot)$  up to  $\epsilon$  accuracy within some safe set of parameters  $S$ , we could exploit the continuity properties to expand the safe set to

$$R_\epsilon(S) := S \cup \bigcap_{i \in \mathcal{I}_g} \{ \mathbf{a} \in \mathcal{A} \mid \exists \mathbf{a}' \in S : g_i(\mathbf{a}') - \epsilon - L \|\mathbf{a} - \mathbf{a}'\| \geq 0 \}, \quad (8)$$

where  $R_\epsilon(S)$  represents the number of parameters that can be classified as safe given that we know  $g$  up to  $\epsilon$ -error inside  $S$  and exploiting the Lipschitz continuity to generalize to new parameters outside of  $S$ . The baseline that we compare against is the limit of repeatedly applying this operator on  $S_0$ ; that is, with  $R_\epsilon^n(S) = R_\epsilon(R_\epsilon^{n-1}(S))$  and  $R_\epsilon^1(S) = R_\epsilon(S)$  the baseline is  $\bar{R}_\epsilon(S_0) := \lim_{n \rightarrow \infty} R_\epsilon^n(S_0)$ . This set contains all the parameters in  $\mathcal{A}$  that could be classified as safe starting from  $S_0$  if we knew the constraint function up to  $\epsilon$  error. This set does not include all the parameters that potentially fulfill the constraints in (1), but is the best we can do without violating the safety constraints. Hence the optimal value that we compare against is not the one in (1), but

<sup>2</sup> The functions  $f$  and  $g_i$  can have different Lipschitz constants  $L_i$ , but we assume a global Lipschitz constant for ease of notation. Additionally, the theoretical results transfer equivalently to the case of Lipschitz-continuity with respect to some metric.

$$f_{\epsilon}^* = \max_{\mathbf{a} \in \mathcal{R}_{\epsilon}(S_0)} f(\mathbf{a}), \quad (9)$$

which is the maximum performance value over the set that we could hope to classify as safe starting from the initial safe set,  $S_0$ .

## 4.1 The algorithm

In this section, we present the SAFEOPT-MC algorithm that guarantees convergence to the previously set baseline. The most critical aspect of the algorithm is safety. However, once safety is ensured, the second challenge is to find an evaluation criterion that enables trading off between exploration, trying to further expand the current estimate of the safe set, and exploitation, trying to improving the estimate of the best parameters within the current set.

To ensure safety, we construct confidence intervals that contain the true functions  $f$  and  $g_i$  with high probability. In particular, we use the posterior GP estimate given the data observed so far. The GP's confidence intervals for the surrogate function in (4) are defined as

$$\mathcal{Q}_n(\mathbf{a}, i) := [\mu_{n-1}(\mathbf{a}, i) \pm \beta_n^{1/2} \sigma_{n-1}(\mathbf{a}, i)], \quad (10)$$

where  $\beta_n$  is a scalar that determines the desired confidence level. This set contains all possible function values between the lower and upper confidence interval based on the GP posterior. The probability of the true function value lying within this interval depends on the choice of  $\beta_n$ , as well as on the properties of the functions that we are modelling with the GP. We provide more details on our assumptions on how to select  $\beta_n$  for our theoretical analysis in Lemma 4.1 and for our experiments in Sect. 4.3.

Rather than defining the lower and upper bounds based on (10), the following analysis requires that consecutive estimates of the lower and upper bounds are contained within each other. This assumption ensures that the safe set does not shrink from one iteration to the next, which we require to prove our results. We relax this assumption in Sect. 4.3. We define the contained set at iteration  $n$  as  $C_n(\mathbf{a}, i) = C_{n-1}(\mathbf{a}, i) \cap \mathcal{Q}_n(\mathbf{a}, i)$ , where  $C_0(\mathbf{a}, i)$  is  $[0, \infty]$  for all  $\mathbf{a} \in S_0$  and  $\mathbb{R}$  otherwise. This ensures that parameters in the initial safe set  $S_0$  remain safe according to the GP model after additional observations. The lower and upper bounds on this set are defined as  $l_n^i(\mathbf{a}) := \min C_n(\mathbf{a}, i)$  and  $u_n^i(\mathbf{a}) := \max C_n(\mathbf{a}, i)$ , respectively. For notational clarity, we write  $l_n^f(\mathbf{a}) := l_n^0(\mathbf{a})$  and  $u_n^f(\mathbf{a}) := u_n^0(\mathbf{a})$  for the performance bounds.

Based on these confidence intervals for the function values and a current safe set  $S_{n-1}$ , we can enlarge the safe set using the Lipschitz continuity properties,

$$S_n = \bigcap_{i \in \mathcal{I}_g} \bigcup_{\mathbf{a} \in S_{n-1}} \{ \mathbf{a}' \in \mathcal{A} \mid l_n^i(\mathbf{a}) - L \|\mathbf{a} - \mathbf{a}'\| \geq 0 \}. \quad (11)$$

The set  $S_n$  contains all points in  $S_{n-1}$ , as well as all additional parameters that fulfill the safety constraints given the GP confidence intervals and the Lipschitz constant.

With the set of safe parameters defined, the last remaining challenge is to trade off between exploration and exploitation. One could, similar to Schreiter et al. (2015), simply select the most uncertain element over the entire set. However, this approach is not sample-efficient, since it involves learning about the entire function rather than restricting evaluations to the relevant parameters. To avoid this, we first define subsets of  $S_n$  that correspond to parameters that could either improve the estimate of the maximum or could expand the safe set. The set of potential maximizers is defined as

$$M_n := \left\{ \mathbf{a} \in S_n \mid u_n^f(\mathbf{a}) \geq \max_{\mathbf{a}' \in S_n} l_n^f(\mathbf{a}') \right\}, \quad (12)$$

which contains all parameters for which the upper bound of the current performance estimate is above the best lower bound. The parameters in  $M_n$  are candidates for the optimum, since they could obtain performance values above the current conservative estimate of the optimal performance.

Similarly, an optimistic set of parameters that could potentially enlarge the safe set is

$$G_n := \{ \mathbf{a} \in S_n \mid e_n(\mathbf{a}) > 0 \}, \quad (13)$$

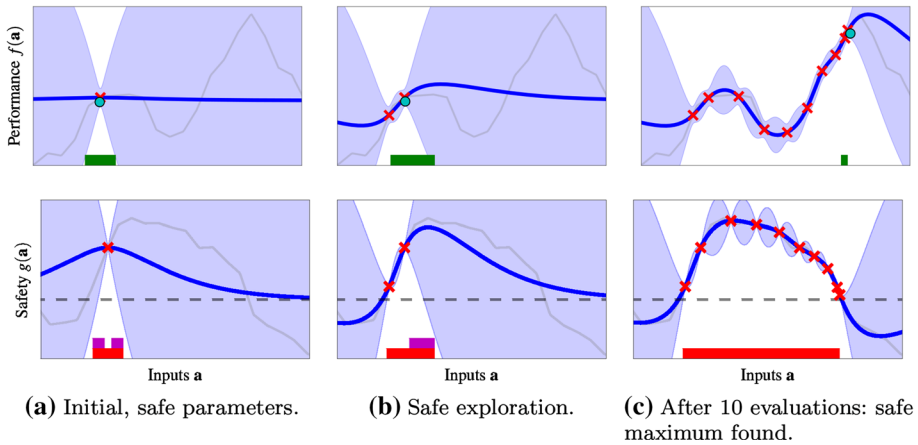
$$e_n(\mathbf{a}) := \left| \{ \mathbf{a}' \in \mathcal{A} \setminus S_n \mid \exists i \in \mathcal{I}_g : u_n^i(\mathbf{a}) - L \|\mathbf{a} - \mathbf{a}'\| \geq 0 \} \right|. \quad (14)$$

The function  $e_n$  enumerates the number of parameters that could additionally be classified as safe if a safety function obtained a measurement equal to its upper confidence bound. Thus, the set  $G_n$  is an optimistic set of parameters that could potentially expand the safe set.

We trade off between the two sets,  $M_n$  and  $G_n$ , by selecting the most uncertain element across all performance and safety functions; that is, at each iteration  $n$  we select

$$\mathbf{a}_n = \operatorname{argmax}_{\mathbf{a} \in G_n \cup M_n} \max_{i \in \mathcal{I}} w_n(\mathbf{a}, i), \quad (15)$$

$$w_n(\mathbf{a}, i) = u_n^i(\mathbf{a}) - l_n^i(\mathbf{a}) \quad (16)$$



**Fig. 2** Optimization with the SAFE OPT-MC algorithm after 1, 2 and 10 parameter evaluations. Based on the mean estimate (blue) and the  $2\sigma$  confidence interval (light blue), the algorithm selects evaluation points for which  $g(\mathbf{a}) \geq 0$  (black dashed) from the safe set  $S_n$  (red), which are either potential maximizers  $M_n$  (green) or expanders  $G_n$  (magenta). It then learns about the function by drawing noisy samples from the unknown, underlying function (light gray). This way, we expand the safe region (red) as much as possible and, simultaneously, find the global optimum of the unknown function (17) (cyan circle) (Color figure online)

as the next parameter set to be evaluated on the real system. The implications of this selection criterion will become more apparent in the next section, but from a high-level view this criterion leads to a behavior that focuses almost exclusively on exploration initially, as the most uncertain points will typically lie on the boundary of the safe set for many commonly used kernels. This changes once the constraint evaluations return results closer to the safety constraints. At this point, the algorithm keeps switching between selecting parameters that are potential maximizers, and parameters that could expand the safe set and lead to new areas in the parameter space with even higher function values. Pseudocode for the algorithm is found in Algorithm 1.

We show an example run of the algorithm in Fig. 2. It starts from an initial safe parameter  $\mathbf{a}_0 \in S_0$  at which we obtain a measurement in Fig. 2a. Based on this, the algorithm uses the continuity properties of the safety function and the GP in order to determine nearby parameters as safe (red set). This corresponds to the region where the high-probability confidence intervals of the GP model (blue shaded) are above the safety threshold (grey dashed line). At the next iteration in Fig. 2b, the algorithm evaluates parameters that are close to the boundary of the safe set, in order to expand the set of safe parameters. Eventually the algorithm converges to the optimal parameters in Fig. 2c, which obtain the largest performance value that is possible without violating the safety constraints. A local optimization approach, e.g. based on estimated gradients<sup>3</sup>, would have gotten stuck in the local optimum at the initial parameter  $\mathbf{a}_0$ .

At any iteration, we can obtain an estimate for the current best parameters from

$$\hat{\mathbf{a}}_n = \operatorname{argmax}_{\mathbf{a} \in S_n} f'_n(\mathbf{a}), \quad (17)$$

which returns the best, safe lower-bound on the performance function  $f$ .

<sup>3</sup> If gradient information is available, it can be incorporated in the GP model too (Solak et al., 2003)

**Algorithm 1: SAFEOPT-MC**


---

**Inputs:** Domain  $\mathcal{A}$ ,  
 GP prior  $k((\mathbf{a}, i), (\mathbf{a}', j))$ ,  
 Lipschitz constant  $L$ ,  
 Initial safe set  $S_0 \subseteq \mathcal{A}$

```

1 for  $n = 1, \dots$  do
2    $S_n \leftarrow \bigcap_{i \in \mathcal{I}_g} \bigcup_{\mathbf{a} \in S_{n-1}} \{\mathbf{a}' \in \mathcal{A} \mid l_n^i(\mathbf{a}) - L\|\mathbf{a} - \mathbf{a}'\| \geq 0\}$ 
3    $M_n \leftarrow \{\mathbf{a} \in S_n \mid u_n^f(\mathbf{a}) \geq \max_{\mathbf{a}' \in S_n} l_n^f(\mathbf{a}')\}$ 
4    $G_n \leftarrow \{\mathbf{a} \in S_n \mid e_n(\mathbf{a}) \geq 0\}$ 
5    $\mathbf{a}_n \leftarrow \operatorname{argmax}_{\mathbf{a} \in G_n \cup M_n} \max_{i \in \mathcal{I}} w_n(\mathbf{a}, i)$ 
6   Measurements:  $\hat{f}(\mathbf{a}_n), \hat{g}_i(\mathbf{a}_n) \forall i = 0, \dots, q$ 
7   Update GP with new data
8 end
```

---

**4.2 Theoretical results**

In this section, we show that the same theoretical framework from the SAFEOPT algorithm (Sui et al., 2015) can be extended to multiple constraints and the evaluation criterion (15). Here, we only provide the results and high-level ideas of the proofs. The mathematical details are provided in Sect. 7. For simplicity, we assume that all function evaluations are corrupted by the same  $\sigma$ -sub-Gaussian noise in this section.

In order to provide guarantees for safety, we need the confidence intervals in (10) to hold for all iterations and functions. In the following, we assume that the surrogate function  $h(\mathbf{a}, i)$  has bounded norm in a reproducing kernel Hilbert space (RKHS, c.f., Christmann and Steinwart (2008)). A RKHS corresponding to a kernel  $k(\cdot, \cdot)$  includes functions of the form  $h(\mathbf{a}, i) = \sum_j \alpha_j k((\mathbf{a}, i), (\mathbf{a}_j, i_j))$  with  $\alpha_i \in \mathbb{R}$  and representer points  $(\mathbf{a}_j, i_j) \in \mathcal{A} \times \mathcal{I}$ . The bounded norm property implies that the coefficients  $\alpha_j$  decay sufficiently fast as  $j$  increases. Intuitively, these functions are well-behaved, in that they are regular with respect to the choice of kernel.

The following Lemma allows us to choose a scaling factor  $\beta_n$  for (10), so that we achieve a specific probability of the true function being contained in the confidence intervals for all iterations.

**Lemma 4.1** (based on Chowdhury and Gopalan (2017)) *Assume that  $h(\mathbf{a}, i)$  has RKHS norm bounded by  $B$  and that measurements are corrupted by  $\sigma$ -sub-Gaussian noise. If  $\beta_n^{1/2} = B + 4\sigma\sqrt{\gamma_{(n-1)|\mathcal{I}|} + 1 + \ln(1/\delta)}$ , then the following holds for all parameters  $\mathbf{a} \in \mathcal{A}$ , function indices  $i \in \mathcal{I}$ , and iterations  $n \geq 1$  jointly with probability at least  $1 - \delta$ :*

$$|h(\mathbf{a}, i) - \mu_{n-1}(\mathbf{a}, i)| \leq \beta_n^{1/2} \sigma_{n-1}(\mathbf{a}, i). \quad (18)$$

Moreover, if the kernel is continuously differentiable, then the corresponding functions are Lipschitz continuous (Christmann & Steinwart, 2008). Note that Lemma 4.1 does not make probabilistic assumptions on  $h$  – in fact,  $h$  could be chosen adversarially, as long as it has bounded norm in the RKHS. Similar results can be obtained for the Bayesian setting where the function  $h$  is assumed to be drawn from the GP prior (Srinivas et al., 2012).

The scaling factor  $\beta_n$  in Lemma 4.1 depends on the *information capacity*  $\gamma_n$  associated with the kernel  $k$ . It is the maximum amount of mutual information that we can obtain about the GP model of  $h(\cdot)$  from  $n$  noisy measurements  $\hat{\mathbf{h}}_{\mathcal{D}}$  at parameters  $\mathcal{D} = \{(\mathbf{a}_1, i_1), \dots\}$ ,

$$\gamma_n = \max_{\mathcal{D} \subseteq \mathcal{A} \times \mathcal{I}, |\mathcal{D}| \leq n} \mathbf{I}(\hat{\mathbf{h}}_{\mathcal{D}}; h). \quad (19)$$

Intuitively, it quantifies a best case scenario where we can select the measurements in the most informative manner possible. The information capacity  $\gamma_n$  has a sublinear dependence on  $n$  for many commonly-used kernels and can numerically approximated up to a small constant factor for any given kernel (Srinivas et al., 2012).

Since the confidence intervals hold with probability  $1 - \delta$  and the safe set is not empty starting from  $S_0$ , it is possible to prove that parameters within the safe set  $S_n$  are always safe with high probability. In order for the algorithm to compete with our baseline, we must additionally ensure that the algorithm learns the true function up to  $\epsilon$  confidence in both the sets  $M_n$  and  $G_n$ . The number of measurements required to achieve this depends on the information capacity  $\gamma_n$ , since it encodes how much information can be obtained about the true function from  $n$  measurements. We use the sublinearity of  $\gamma_n$  in order to bound the number of samples required to estimate the function up to  $\epsilon$  accuracy. We have the following result:

**Theorem 4.1** Assume that  $h(\mathbf{a}, i)$  has bounded norm in an RKHS and that the measurement noise is  $\sigma$ -sub-Gaussian. Also, assume that  $S_0 \neq \emptyset$  and  $g_i(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in S_0$  and  $i \in \mathcal{I}_g$ . Choose  $\beta_n$  as in Lemma 4.1, define  $\hat{\mathbf{a}}_n$  as in (17), and let  $n^*(\epsilon, \delta)$  be the smallest positive integer satisfying

$$\frac{n^*}{\beta_{n^*} \gamma_{|\mathcal{I}|n^*}} \geq \frac{\tilde{c}(|\bar{R}_0(S_0)| + 1)}{\epsilon^2}, \quad (20)$$

where  $\tilde{c} = 8 / \log(1 + \sigma^{-2})$ . For any  $\epsilon > 0$  and  $\delta \in (0, 1)$ , when running Algorithm 1 the following inequalities jointly hold with probability at least  $1 - \delta$ :

1. *Safety*:  $\forall n \geq 1, \forall i \in \mathcal{I}_g : g_i(\mathbf{a}_n) \geq 0$
2. *Optimality*:  $\forall n \geq n^*, f(\hat{\mathbf{a}}_n) \geq f_{\epsilon}^* - \epsilon$

**Proof** Main idea: safety follows from Lemma 4.1, since accurate confidence intervals imply that we do not evaluate unsafe parameters. For the optimality, the main idea is that, since we evaluate the most uncertain element in  $M_n \cup G_n$ , the uncertainty about the maximum is bounded by  $w_n(\mathbf{a}_n, i_n)$ . The result follows from showing that, after a finite number of evaluations, either the safe set expands or the maximum uncertainty within  $M_n \cup G_n$  shrinks to  $\epsilon$ . See Sect. 7 for derivations and details.  $\square$

Theorem 4.1 states that, given the assumptions we made about the underlying function, Algorithm 1 explores the state space without violating the safety constraints and, after at most  $n^*$  samples, finds an estimate that is  $\epsilon$ -close to the optimal value over the safely reachable region. The information capacity  $\gamma_{|\mathcal{I}|n^*}$ , grows at a faster rate of  $|\mathcal{I}|n$  compared to  $n$  in SAFEOPT, since we obtain  $|\mathcal{I}|$  observations at the same parameters set  $\mathbf{a}$ ,

while the  $\text{SAFEOPT}$  analysis assumes every measurement is optimized independently. However,  $\gamma_{|\mathcal{Z}|n}$  remains sublinear in  $n$ , see Sect. 7.

#### 4.2.1 Contexts

In this section, we show how the theoretical guarantees derived in the previous section transfer to contextual Bayesian optimization. In this setting, part of the variables that influence the performance, the contexts, are fixed by an external process that we do not necessarily control. Since the contexts are part of the kernel function in (6), we make explicit assumptions about how smoothly the functions vary with the context. In normal Bayesian optimization, it was shown by Krause and Ong (2011) that an algorithm that optimizes the GP-UCB criterion in (7) for a fixed context converges to the global optimum after repeatedly seeing the corresponding context.

Intuitively, the fact that part of the variables that influence the performance, the contexts, are now fixed by an external process should not impact the algorithm on a fundamental level. However, safety is a critical issue in our experiments and, in general, one could always select an adversarial context for which we do not have sufficient knowledge to determine safe controller parameters. As a consequence, we make the additional assumption that only ‘safe’ contexts are visited; that is, we assume the following:

**Assumption 1** For any  $n \geq 1$ , the context  $\mathbf{z}_n \in \mathcal{Z}$  is selected such that  $S_n(\mathbf{z}_n) \neq \emptyset$ .

Here,  $S_n(\mathbf{z}_n)$  denotes the safe set for the given context  $\mathbf{z}_n$ . Intuitively, Assumption 1 ensures that for every context there exists at least one parameter choice that is known to satisfy all safety constraints. This assumption includes the case where safe initial parameters for all contexts are known *a priori* and the case where the algorithm terminates and asks for help from a domain-expert whenever a context leads to an empty safe set.

A trivial extension of  $\text{SAFEOPT-MC}$  to contexts is to run  $|\mathcal{Z}|$  independent instances of Algorithm 1, one for each context. This way, it is sufficient to repeatedly see a context several times to apply the previous results to the safe contextual optimization case. One can apply the previous analysis to this setting, but it would fail to yield guarantees that hold *jointly* for all contexts.

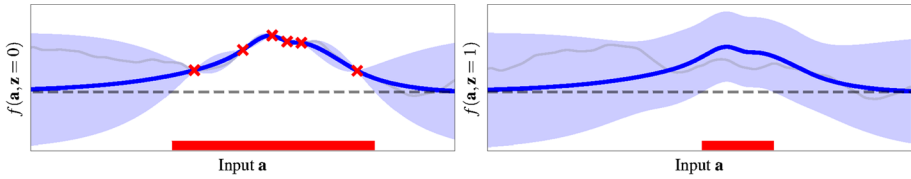
In order to obtain stronger results that hold jointly across all contexts in  $\mathcal{Z}$ , we adapt the information capacity (worst-case mutual information)  $\gamma_n$  to consider contexts,

$$\gamma_n = \max_{\mathcal{D} \subseteq \mathcal{A} \times \mathcal{Z} \times \mathcal{I}_n, |\mathcal{D}| \leq n} \mathbf{I}(\hat{\mathbf{h}}_{\mathcal{D}}; h), \quad (21)$$

Unlike in (19), the mutual information is maximized across contexts in (21). As a result, we can use Lemma 4.1 to obtain confidence intervals that hold jointly across all contexts.

A second challenge is that contexts are chosen in an arbitrary order. This is in stark contrast to the parameters  $\mathbf{a}_n$ , which are chosen according to (15) in order to be informative. This means that any tight finite sample bound on Algorithm 1 must necessarily depend on the order of contexts. The following theorem accounts for both of these challenges.

**Theorem 4.2** *Under the assumptions of Theorem 4.1 and Assumption 1. Choose  $\beta_n$  as in Lemma 4.1, where  $\gamma_n$  is now the worst-case mutual information over contexts as in (21). For a given context order  $(\mathbf{z}_1, \mathbf{z}_2, \dots)$  and any context  $\mathbf{z} \in \mathcal{Z}$ , let*



**Fig. 3** Example run of the context-dependent SAFEOPT-MC algorithm. For the first six samples, the algorithm only sees the context  $\mathbf{z} = 0$  (left function) and obtains measurements there (red crosses). However, by exploiting correlations between different contexts, the algorithm can transfer knowledge about the shape of the function and safe set over to a different context,  $\mathbf{z} = 1$  (right function). This enables the algorithm to be significantly more data-efficient (Color figure online)

$$n(\mathbf{z}) = \sum_{n=1}^{n^*(\mathbf{z})} \mathbb{1}_{\mathbf{z}=\mathbf{z}_n} \quad (22)$$

be the number of times that we have observed the context  $\mathbf{z}$  up to iteration  $n^*$  and  $\mathbb{1}$  is the indicator function. Let  $n^*(\mathbf{z})$  be the smallest positive integers such that

$$\frac{n(\mathbf{z})}{\beta_{n^*(\mathbf{z})} \gamma_{n(\mathbf{z})|\mathcal{I}}(\mathbf{z})} \geq \frac{\tilde{c}(|\bar{R}_0(S_0(\mathbf{z}))| + 1)}{\epsilon^2}, \quad (23)$$

where  $\tilde{c} = 8/\log(1 + \sigma^{-2})$ . We denote the information capacity for a fixed context  $\mathbf{z}$  by  $\gamma_n(\mathbf{z})$ . That is, with  $h_{\mathbf{z}}(\mathbf{a}, i) = h(\mathbf{a}, i, \mathbf{z})$  it is defined as  $\gamma_n(\mathbf{z}) = \max_{\mathcal{D} \subseteq \mathcal{A} \times \mathcal{I}, |\mathcal{D}| \leq n} \mathbb{I}(\hat{\mathbf{h}}_{\mathbf{z}}|_{\mathcal{D}}; h_{\mathbf{z}})$ . For any  $\epsilon > 0$  and  $\delta \in (0, 1)$ , let  $f_{\epsilon}^*(\mathbf{z}) = \max_{\mathbf{a} \in \bar{R}_{\epsilon}(S_0)} f(\mathbf{a}, \mathbf{z})$ . Then, when running Algorithm 1 the following inequalities jointly hold with probability at least  $1 - \delta$ :

1.  $\forall n \geq 1, i \in \mathcal{I}_g : g_i(\mathbf{a}_n, \mathbf{z}_n) \geq 0$
2.  $\forall \mathbf{z} \in \mathcal{Z}, n \geq n^*(\mathbf{z}) : f(\hat{\mathbf{a}}_n, \mathbf{z}) \geq f_{\epsilon}^*(\mathbf{z}) - \epsilon$

**Proof** For a fixed context,  $\mathbf{z}_n = \mathbf{z} \forall n$ , we have  $n^*(\mathbf{z}) = n(\mathbf{z})$  and the results follow directly as in Theorem 4.1. Otherwise, the only difference in the proofs is that  $\beta$  depends on the information capacity over contexts, making sure that the confidence intervals are valid across contexts. By visiting contexts in  $\mathcal{Z} \setminus \{\mathbf{z}\}$ , we obtain more measurements and  $\beta$  increases, which in turn increases the upper bound on the number of samples required at context  $\mathbf{z}$ .  $\square$

Theorem 4.2 states that the contextual variant of Algorithm 1 enjoys the same safety guarantees as the non-contextual version. Additionally, it shows that, after gaining enough information about a particular context, it can identify the optimal parameters. In practice, this upper bound is conservative, since it does not account for knowledge transfer across contexts. In practice, correlations between contexts significantly speed up the learning process. For example, in Fig. 3 we show a contextual safe optimization problem with two contexts. Even though the algorithm has only been able to explore the parameter space at the first context ( $\mathbf{z} = 0$ , left function), the correlation between the functions means that information can be transferred to the as-of-yet unobserved context ( $\mathbf{z} = 1$ , right function).



This knowledge transfer significantly improves data-efficiency and the number of evaluations required by the algorithm.

### 4.3 Practical implementation

In this section, we discuss possible changes to Algorithm 1 that make the algorithm more practical, at the expense of loosing some of the theoretical guarantees. One challenge in applying Algorithm 1 in practice, is defining a suitable Lipschitz constant. In particular, specifying the wrong constant can lead to conservativeness or unsafe parameters being evaluated. Moreover, smoothness assumptions are already encoded by the kernel choice, which is more intuitive to specify than Lipschitz constants on their own. In practice, we use only the GP model to ensure safety (Berkenkamp et al., 2016b); that is, we define  $l_n^i(\mathbf{a}) = \min Q_n(\mathbf{a}, i)$  and  $u_n^i(\mathbf{a}, i) = \max Q_n(\mathbf{a}, i)$  in terms of the confidence intervals of the GP directly. Thus, we can define the safe set without a Lipschitz constant as

$$S_n = S_0 \cup \{\mathbf{a} \in \mathcal{A} \mid \forall i \in \mathcal{I}_g : l_n^i(\mathbf{a}) \geq 0\}. \quad (24)$$

While it is difficult to prove the full exploration of the safely reachable set as in Theorem 4.1, the resulting algorithm remains safe:

**Lemma 4.2** *With the assumptions of Lemma 4.1,  $S_0 \neq \emptyset$ , and  $g_i(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in S_0$  and  $i \in \mathcal{I}_g$ , when running Algorithm 1 with the safe set defined as in (24), the following holds with probability at least  $1 - \delta$ :*

$$\forall n \geq 1, \forall i \in \mathcal{I}_g : g_i(\mathbf{a}_n) \geq 0. \quad (25)$$

**Proof** The confidence intervals hold with probability  $1 - \delta$  following Lemma 4.1. Since  $S_n$  in (24) is defined as the set of parameters that fulfill the safety constraint and the safe set is never empty since  $S_0 \neq \emptyset$ , the claim follows.  $\square$

Similarly, the set of expanders can be defined in terms of the GP directly, by adding optimistic measurements and counting the number of new parameters that are classified as safe, see Berkenkamp et al. (2016b) for details. However, this potentially adds a large computational burden.

The parameter  $\beta_n$ , which determines the GP's confidence interval in Lemma 4.1, may be impractically conservative for experiments. The theoretical safety results also hold when we replace  $\gamma_n$  in  $\beta_n$  by the empirical mutual information gained so far,  $I(\hat{\mathbf{h}}_{\mathcal{D}_n \times \mathcal{I}}, h)$ . Empirically, depending on the application, one may also consider setting  $\beta_n$  to a constant value. This roughly corresponds to bounding the failure probability per iteration, rather than over all iterations.

Learning all the different functions,  $f$  and  $g_i$ , up to the same accuracy  $\epsilon$  may be restrictive if they are scaled differently. A possible solution is to either scale the observed data, or to scale the uncertainties in (15) by the prior variances of the kernels for that specific output. This enables (15) to make more homogeneous decisions across different scales.

## 5 Experiments

In this section, we demonstrate Algorithm 1 (with the changes discussed in Sect. 4.3) in experiments on a quadrotor vehicle, a Parrot AR.Drone 2.0.

A Python implementation of the SAFEPT-MC algorithm that builds on GPy (2012), a GP library, is available at <http://github.com/befelix/SafeOpt>. Videos of the experiments can be found at

- Section 5.3: [http://tiny.cc/icra16\\_video](http://tiny.cc/icra16_video)
- Section 5.4: <https://youtu.be/rLmwYtoE3yg>
- Section 5.5: <https://youtu.be/4xC4OSiIDGk>

### 5.1 Experimental setup

During the experiments, measurements of all vehicle states are estimated from position and pose data provided by an overhead motion capture camera system. The quadrotor's dynamics can be described by 12 states: positions  $\mathbf{x} = (x, y, z)$ , velocities  $\dot{\mathbf{x}} = (\dot{x}, \dot{y}, \dot{z})$ , ZYX Euler angles  $(\phi, \theta, \psi)$ , and body angular velocities  $(\omega_x, \omega_y, \omega_z)$ . The control inputs  $\mathbf{u}$  are the desired roll and pitch angles  $\theta_{\text{des}}$  and  $\phi_{\text{des}}$ , the desired  $z$ -velocity  $\dot{z}_{\text{des}}$ , and the desired yaw angular velocity  $\omega_{z,\text{des}}$ , which in turn are inputs to an unknown, proprietary, on-board controller.

The position dynamics in the global coordinate frame are

$$\ddot{\mathbf{x}} = R_{ZYX}(\phi, \theta, \psi)\mathbf{f} - \mathbf{g}, \quad (26)$$

where  $R_{ZYX}$  is the rotation matrix from the body frame to the inertial frame,  $\mathbf{f} = (0, 0, c)$  is the mass-normalized thrust, and  $\mathbf{g} = (0, 0, g)$  is the gravitational force. The goal of the controller is to track a reference signal. We assume that  $z$ -position and the yaw angle are controlled by fixed control laws and focus on the position control in  $x$ - and  $y$  direction. We use two different control laws in the following experiments.

The most simple control law that can be used for this setting is a PD-controller, defined as

$$\phi_{\text{des}} = k_1(x_k - x_{\text{des}}) + k_2(\dot{x} - \dot{x}_{\text{des}}), \quad (27)$$

$$\theta_{\text{des}} = k_1(y_k - y_{\text{des}}) + k_2(\dot{y} - \dot{y}_{\text{des}}), \quad (28)$$

where  $\mathbf{a} = (k_1, k_2)$  are the two tuning parameters. Intuitively, a larger parameter  $k_1$  encourages tracking reference changes more aggressively, while  $k_2$  is a damping factor that encourages moderate velocities.

A more sophisticated approach to control quadrotor vehicles is to use estimates of the angles and accelerations to solve for the thrust  $c$ . We use loop shaping on the horizontal position control loops so that they behave in the manner of a second-order systems with time constant  $\tau$  and damping ratio  $\zeta$ . Based on a given desired reference trajectory, commanded accelerations are computed as

$$\ddot{x}_c = \frac{1}{\tau^2}(x_{\text{des}} - x) + \frac{2\zeta}{\tau}(\dot{x}_{\text{des}} - \dot{x}), \quad (29)$$

$$\ddot{y}_c = \frac{1}{\tau^2}(y_{\text{des}} - y) + \frac{2\zeta}{\tau}(\dot{y}_{\text{des}} - \dot{y}). \quad (30)$$

From the commanded accelerations, we then obtain the control inputs for the desired roll and pitch angles by solving (26) for the angles. Here, the tuning parameters are  $\mathbf{a} = (\tau, \zeta)$ . For details regarding the controllers see Schoellig et al. (2012); Lupashin et al. (2014).

The quadrotor was controlled using the `ardrone_autonomy` and `vicon_bridge` packages in ROS Hydro. Computations for the `SAFE`OPT-MC algorithm in Algorithm 1 were conducted on a regular laptop and took significantly less than one second per iteration. As a result, experiments could be conducted continuously without interruptions or human interventions.

## 5.2 Kernel selection

Algorithm 1 critically depends on the GP model for the performance and constraint functions. In this section, we review the kernel used in our experiments and the kind of models that they encode. A more thorough review of kernel properties can be found in Rasmussen and Williams (2006).

In our experiments, we use the Matérn kernel with parameter  $\nu = 3/2$  (Rasmussen & Williams, 2006),

$$k(\mathbf{a}, \mathbf{a}') = \kappa^2 \left(1 + \sqrt{3} r(\mathbf{a}, \mathbf{a}')\right) \exp\left(-\sqrt{3} r(\mathbf{a}, \mathbf{a}')\right), \quad (31)$$

$$r(\mathbf{a}, \mathbf{a}') = \sqrt{(\mathbf{a} - \mathbf{a}')^T \mathbf{M}^{-2} (\mathbf{a} - \mathbf{a}')}, \quad (32)$$

which encodes that mean functions are one-times differentiable. This is in contrast to the commonly used squared exponential kernels, which encode smooth (infinitely differentiable) functions. With the Matérn kernel, the GP model is parameterized by three hyperparameters: measurement noise  $\sigma^2$  in (2) and (3), the kernel's prior variance  $\kappa^2$ , and positive lengthscales  $\mathbf{l} \in \mathcal{R}_+^A$ , which are the diagonal elements of the diagonal matrix  $\mathbf{M}$ ,  $\mathbf{M} = \text{diag}(\mathbf{l})$ . These hyperparameters have intuitive interpretations: the variance of the measurement noise  $\kappa^2$  corresponds to the noise in the observations, which includes any randomness in the algorithm and initial conditions, and random disturbances. The prior variance  $\kappa^2$  determines the expected magnitude of function values; that is,  $|f(\mathbf{a})| \leq \kappa$  with probability 0.68 according to the GP prior. Lastly, the lengthscales  $\mathbf{l}$  determine how quickly the covariance between neighboring values deteriorates with their distance. The smaller the lengthscales, the faster the function values can change from one parameter set to the next. In particular, the high-probability Lipschitz constant encoded by this kernel depends on the ratio between the prior variance and the lengthscales,  $\kappa/\mathbf{l}$ .

When using GPs to model dynamic systems, typically a maximum likelihood estimate of the hyperparameters is used based on data; see Ostafew et al. (2016) for an example. For Bayesian optimization, the GP model is used to actively acquire data, rather than only using it for regression based on existing data. This dependence between the kernel hyperparameters and the acquired data is known to lead to poor results in Bayesian optimization when using a maximum likelihood estimate of the hyperparameters during data acquisition (Bull, 2011). In particular, the corresponding GP estimate is not guaranteed to contain the true function as in Lemma 4.1. In this work, we critically rely on these confidence bounds to guarantee safety. As a consequence, we do not adapt the hyperparameters

as more data becomes available, but treat the kernel as a prior over functions in the true Bayesian sense; that is, the kernel hyperparameters encode our prior knowledge about the functions that we model and are fixed before experiments begin. While this requires intuition about the process, if we chose the parameters more conservatively than required (e.g., short lengthscales), the safety guarantees still hold (c.f., Bull (2011)) but the optimization becomes less data-efficient. As such, safety can be ensured by picking a conservative model at the cost of optimization efficiency.

### 5.3 Linear control

In this section, we use SAFEOPT-MC to optimize the parameters of the linear control law in (27). The entire control algorithm consists of this control law together with the on-board controller and state estimation.

The goal is to find controller parameters that maximize the performance during a 1-meter reference position change. For an experiment with parameters  $\mathbf{a}_n$  at iteration  $n$ , the performance function is defined as

$$f(\mathbf{a}_n) = c(\mathbf{a}_n) - 0.95 c(\mathbf{a}_0), \quad (33)$$

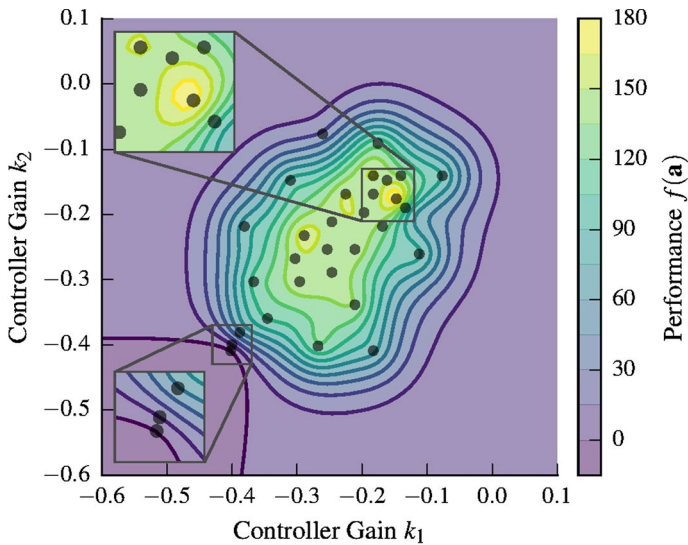
$$c(\mathbf{a}_n) = - \sum_{k=0}^N \mathbf{x}_k^T \mathbf{Q} \mathbf{x}_k + R u_k^2, \quad (34)$$

where, to compute the cost  $c$ , the states  $\mathbf{x} = (x - 1, \dot{x}, \phi, \omega)$  and the input  $u$  are weighted by positive semi-definite matrices  $\mathbf{Q}$  and  $R$ . The subscript  $k$  indicates the state measurement at time step  $k$  in the trajectory and the time horizon is 5 s ( $N = 350$ ). Performance is defined as the cost improvement relative to 95% of the initial controller cost. The safety constraint is defined only in terms of the performance; that is, the constraint is  $g(\mathbf{a}) = f(\mathbf{a}) \geq 0$ , which encodes that we do not want to evaluate controller parameters that perform significantly worse than the initial parameters.

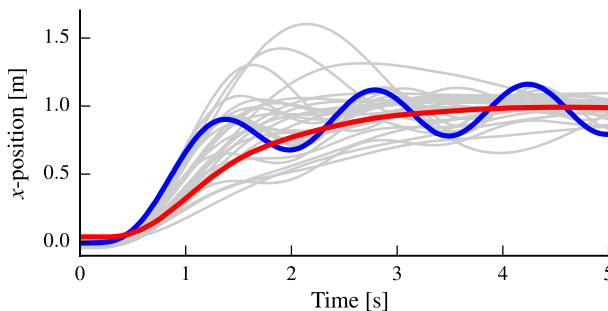
While the optimal controller parameters could be easily computed given an accurate model of the system, we do not have a model of the dynamics of the proprietary, on-board controller and the time delays in the system. Moreover, we want to optimize the performance for the real, nonlinear quadrotor system, which is difficult to model accurately. An inaccurate model of the system could be used to improve the prior GP model of the performance function, with the goal of achieving faster convergence. In this case, the uncertainty in the GP model of the performance function would account for inaccuracies in the system model.

We define the domain of the controller parameters as  $[-0.6, 0.1]^2$ , explicitly including positive controller parameters that certainly lead to crashes. In practice, one would exclude parameters that are known to be unsafe *a priori*. The initial controller parameters are  $(-0.4, -0.4)$ , which result in a controller with poor performance. Decreasing the controller gains further leads to unstable controllers.

The parameters for the experiments were set as follows: the length-scales were set to 0.05 for both parameters, which corresponds to the notion that a 0.05–0.1 change in the parameters leads to very different performance values. The prior standard deviation,  $\kappa$ , and the noise standard deviation,  $\sigma$ , are set to 5% and 10% of the performance of the initial controller,  $f(\mathbf{a}_0)$ , respectively. The noise standard deviation,  $\sigma$ , mostly models errors due to initial position offsets, since state measurements have low noise. The size of these errors



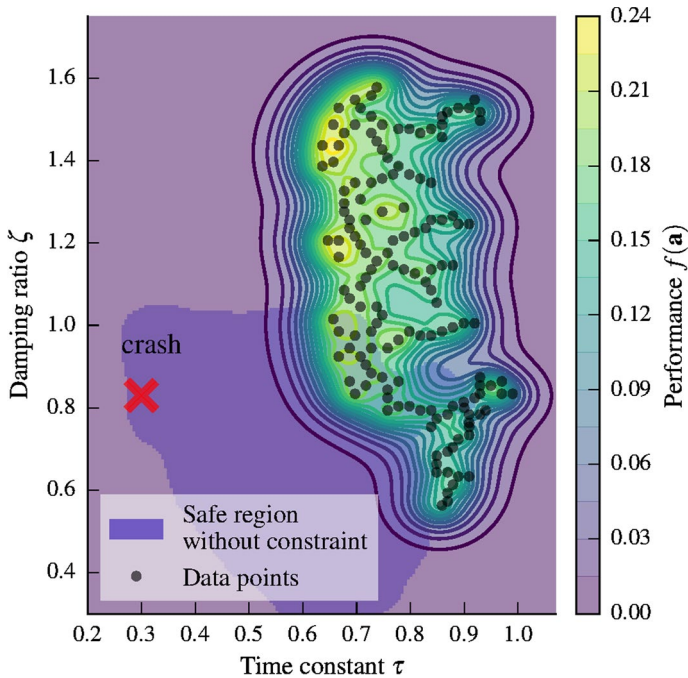
**Fig. 4** GP mean estimate of the performance function after 30 evaluations. The algorithm adaptively decides which parameters to evaluate based on safety and informativeness. In the bottom-left corner, there is the magnified section of the first three samples, which are close together to determine the location of the initial, safe region. The maximum, magnified in the top-left corner, also has more samples to determine the precise location of the maximum. Other areas are more coarsely sampled to expand the safe region



**Fig. 5** The quadrotor controller performance is evaluated during a 5 s evaluation interval, where a 1 m reference position change must be performed. The trajectories correspond to the optimization routine in Fig. 4. The initial controller (blue) performs poorly but is stable. In contrast, the optimized controller (red) shows an optimized, smooth, and fast response. The trajectories of other controller parameters that were evaluated are shown in gray (Color figure online)

depends on the choice of the matrices  $\mathbf{Q}$  and  $\mathbf{R}$ . By choosing  $\sigma$  as a function of the initial performance, we account for the  $\mathbf{Q}$  and  $\mathbf{R}$  dependency. Similarly,  $\kappa$  specifies the expected size of the performance function values. Initially, the best we can do is to set this quantity dependent on the initial performance and leave additional room for future, larger performance values. For the GP model, we choose  $\beta_n^{1/2} = 2$  to define the confidence interval in (10).

The resulting, estimated performance function after running Algorithm 1 for 30 experiments is shown in Fig. 4. The unknown function has been reliably identified. Samples are



**Fig. 6** Mean estimate of the root-mean-square error when optimizing the parameters of the nonlinear control law for a step response, subject to safety constraints. The algorithm carefully evaluates only safe parameter combinations, until the safe region cannot be expanded further without violating constraints. Without the safety constraint, the algorithm explores a larger region of the parameter space (light blue) and eventually evaluates an unsafe parameter set (Color figure online)

spread out over the entire safe set, with more samples close to the maximum of the function and close to the initial controller parameters. No unsafe parameters below the safety threshold were evaluated on the real system.

Typically, the optimization behavior of Algorithm 1 can be roughly separated into three stages. Initially, the algorithm evaluates controller parameters close to the initial parameters in order for the GP to acquire information about the safe set (see lower-left, zoomed-in section in Fig. 4). Once a region of safe controller parameters is determined, the algorithm evaluates the performance function more coarsely in order to expand the safe set. Eventually, the controller parameters are refined by evaluating high-performance parameters that are potential maximizers in a finer grid (see upper-left, zoomed-in section in Fig. 4). The trajectories of the initial, best and intermediate controllers can be seen in Fig. 5.

## 5.4 Nonlinear control

In the previous section, we showed how to optimize the performance of a linear control law subject to a simple constraint on performance. In this section, we optimize the nonlinear controller in (29) and (30) and show how more complex constraints can be used.

We use the same task as in the previous section, but this time the goal is to minimize the root-mean-square error (RMSE) over a time horizon of 5 s ( $N = 350$  samples) during a 1-meter reference position change in  $x$ -direction. We define the performance function,

$$f(\mathbf{a}_n) = c(\mathbf{a}_n) - 0.75 c(\mathbf{a}_0), \quad (35)$$

$$c(\mathbf{a}_n) = \frac{1}{\sqrt{N}} \left( \sum_{k=1}^N \|\mathbf{x}_k - \mathbf{x}_{\text{des},k}\|_2^2 \right)^{1/2}, \quad (36)$$

as the performance relative to 75% of the performance of the initial parameters  $\mathbf{a}_0 = (0.9, 0.8)$ . We define the GP model of this performance function as follows: in this experiment, measurement noise is minimal, since the positions are measured accurately by the overhead camera system. However, to capture errors in the initial position, we define  $\sigma = 0.05c(\mathbf{a}_0)$ . We assume that we can improve the initial controller by roughly 20%, so we set  $\kappa = 0.2c(\mathbf{a}_0)$ . The lengthscales are set to 0.05 in order to encourage cautious exploration. These parameters turned out to be conservative for the real system. Notice that the cost is specified relative to  $c(\mathbf{a}_0)$  instead of  $f(\mathbf{a}_0)$  as in Sect. 5.3. Since  $c(\mathbf{a}_0) > f(\mathbf{a}_0)$ , these hyperparameters are more conservative, so that we require more evaluations on the real system. The reason for this more conservative choice is that the nonlinear controller is expected to have a less smooth performance function, unlike the one for linear control, which is expected to be roughly quadratic.

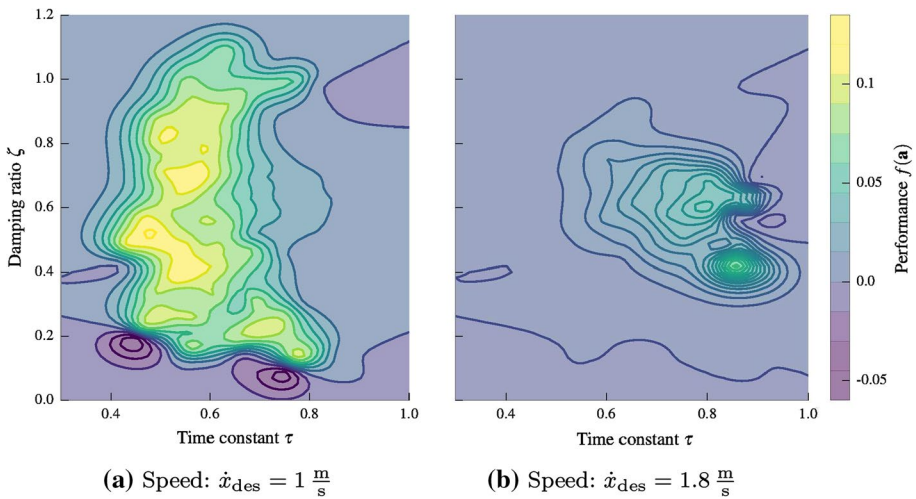
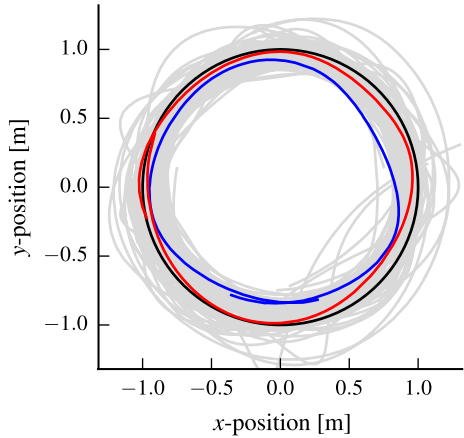
If, as in the previous section, one were to set the safety constraint to  $g_1(\mathbf{a}) = f(\mathbf{a})$ , the algorithm would classify the blue shaded region in Fig. 6 as safe. This region includes time constants as low as  $\tau = 0.3$ , which encourage highly aggressive maneuvers, as would be expected from a performance function that encourages changing position as fast as possible. However, these high gains amplify noise in the measurements, which can lead to crashes; that is, the performance-based constraint cannot properly encode safety. Notice that the blue shaded area does not correspond to full exploration, since the experiment was aborted after the first, serious crash. The reason for the unsafe exploration is that the RMSE performance function in (36) does not implicitly encode safety the same way as weighting of state and input errors in Fig. 4 does, where high control gains are penalized due to the control cost. In particular, this means that if the gains are chosen too high so that the robot crashes, we have an abrupt drop in performance that violates the theoretical assumptions that allow us to use a GP model for exploration. To avoid this, we need to specify additional safety constraints.

One indication of unsafe behavior in quadrotors are high angular velocities when the quadrotor oscillates around the reference point. We define an additional safety constraint on the maximum angular velocity  $\max_k |\omega_{x,k}| \leq 0.5$  rad/s by setting  $g_2(\mathbf{a}) = 0.5 - \max_k |\omega_{x,k}|$ . The corresponding hyperparameters are selected as  $\sigma_2 = 0.1$ ,  $l = 0.2$ , and  $\kappa = 0.25$ . The measurement noise can be chosen to be relatively small here, since it corresponds to a single measurement of angular velocity. Note that it is difficult to perform the step maneuver with an angular velocity lower than 0.4 rad/s, so that typical values of  $g_2$  are smaller than 0.1.

With this additional safety constraint, the algorithm explores the parameter space and stops before the safety constraints are violated, as can be seen in Fig. 6. Rather than exploring smaller time constants  $\tau$  (higher gains), the algorithm evaluates larger damping ratios, which allow slightly smaller values of  $\tau$  and therefore higher performance without violating the safety constraints. The optimal parameters are to the top-left of the safe set, where small time constants encourage tracking the reference aggressively, while the increased damping ratio ensures a moderate angular velocity.



**Fig. 7** The trajectories (gray) resulting from iteratively optimizing the controller parameters for a unit circle reference trajectory at 1 m/s (black). The trajectory with the initial parameters (blue) has poor tracking performance, while the optimized parameters (red) perform significantly better. The flight is safe, i.e., only safe parameters are evaluated (Color figure online)



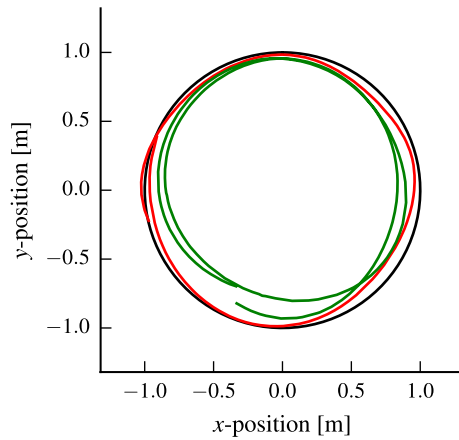
**Fig. 8** The mean estimate of the performance function for the circle trajectory in Fig. 7 for a speed of 1 m/s (left) and 1.8 m/s (right). Extending the kernel with a context for speed allows to transfer knowledge to different speeds and leads to speed-dependent optimal control parameters, speeding up the learning for higher speeds

## 5.5 Circle trajectory

In this section, we use the same nonlinear controller and cost function as in the previous section, but aim to optimize the RMSE with respect to a circle trajectory of radius 1 m at a speed of 1 m/s. The reference is defined as a point moving along the circle at the desired speed. Feasibility of such motions has been analyzed in Schoellig et al. (2011).

In order to ensure good tracking behavior, we define safety as a constraint on the maximum RMSE of 0.2 m. Additionally, we use the same constraint on the maximum angular velocity around the  $x$  and  $y$  axis of 0.5 rad/s as before. The yaw-angle is set so that the quadrotor always points to the center of the circle, which ideally should lead to zero





**Fig. 9** Trajectories with optimal parameters for speeds of 1 m/s (red) and 1.8 [m/s] (green) when tracking the black reference. At slower speeds there exist aggressive controller parameters that allow the quadrotor to track the reference almost perfectly. At higher speeds, actuator saturation limits the achievable performance. Due to the safe optimization framework the maximum speed can be found that does not deviate more from the reference trajectory than is allowed by the safety constraint. The corresponding performance functions can be seen in Fig. 8 (Color figure online)

angular velocity. Deviations from this are an indication of unsafe behavior. We use the same kernel hyperparameters as in Sect. 5.4.

The trajectories that result from running the optimization algorithm are shown in Fig. 7. The initial controller parameters lead to very poor performance. In particular, the initial time constant is too large, so that the quadrotor lags behind the reference. As a result, the quadrotor flies a circle of smaller radius. In contrast, the resulting optimized trajectory (in red) is the best that can be obtained given the safety constraints and controller structure above. The mean estimate of the performance function after the experiments can be seen in Fig. 8a. The optimal parameters have smaller time constants, so that the position is tracked aggressively. Since the reference point moves of 1 m/s, these aggressive controller parameters do not lead to unsafe behavior. During the entire optimization, only safe parameters that keep the vehicle within the constraints on RMSE and angular velocity are evaluated.

## 5.6 Context-dependent optimization

In this section, we show how the knowledge about good controller parameters at low speeds can be used to speed up the safe learning at higher speeds.

In our circle experiment, the quadrotor tracked a moving reference. As this reference moves with high velocities, the quadrotor gets pushed to its physical actuator limits and starts to lag behind the reference. This causes the circle that is flown by the quadrotor to have a smaller radius than the reference trajectory. In this section, the goal is to maximize the speed of the reference trajectory subject to the safety constraints of the previous experiment in Sect. 5.5. One way to achieve this goal, is to add the speed of the reference point to the performance function. However, this would lead to more experiments, as the algorithm will explore the safe parameter space for every velocity. Instead, here we define the trajectory speed as a context, which is set externally. In particular, we set

$$z_n = \underset{v \in \mathbb{R}, \mathbf{a} \in \mathcal{A}}{\operatorname{argmax}} v \quad \text{subject to: } g_i(\mathbf{a}, v) \geq 0, \forall i \in \mathcal{I}_g, \quad (37)$$

that is, we select the maximum velocity for which there are safe parameters known. While here we select the context manually, in practice contexts can be used to model any external, measurable variables, such as the battery level, see Sect. 3.3.

In order to transfer knowledge about good controller parameters from the slow speed in Sect. 5.5 to higher speeds, we model how performance and constraints vary with desired speed by defining a kernel  $k_z(\dot{x}_{\text{des}}, \dot{x}'_{\text{des}})$  over contexts. We use the same kernel structure as in (6) and hyperparameters  $\kappa = 1$  and  $l = 0.25$ . Based on the data from Sect. 5.5, the extended model allows us to determine speeds for which safe controller parameters are known.

Starting from the data of the previous experiments in Sect. 5.5, we run SAFEOPT-MC using the extended kernel with the additional speed context determined by (37). This allows us to find optimal parameters for increasingly higher speeds, which satisfy the constraints. We can safely increase the speed up to 1.8 m/s. We show the mean performance function estimates for two speeds in Fig. 8. For lower speeds, the best controller parameters track the reference position more aggressively (low  $\tau$ ). For higher speeds, this behavior becomes unsafe as the quadrotor lags behind the reference point. Instead, the optimal parameters shift to higher time constants (lower gains). Additionally, as expected, high speeds lead to higher reference tracking errors. Increasing the reference velocity any further causes the performance constraint to be violated.

The trajectories that result from applying the optimal parameters for a speed of 1 m/s and the maximum safe speed of 1.8 m/s can be seen in Fig. 9. For the relatively slow speed of 1 m/s the quadrotor can track the circle well using aggressive parameters. For the higher speed, the reference trajectory moves too fast for the quadrotor to track perfectly within the actuator limits, so that the best parameters just barely satisfy the safety constraint on the average deviation from the reference. Overall, this approach allows us to find context-dependent parameters, while remaining within the safety constraints.

## 6 Conclusion and future work

We presented a generalization of the Safe Bayesian Optimization algorithm of Sui et al. (2015) that allows multiple, separate safety constraints to be specified and applied it to non-linear control problems on a quadrotor vehicle. Overall, the algorithm enabled efficient and automatic optimization of parameters without violating the safety constraints, which would lead to system failures. Currently, the algorithm is mostly applicable to low-dimensional problems due to the computational burden of optimizing (15) and the statistical problem of defining suitable GP priors in high-dimensions. While interesting progress has been made in this direction in the standard Bayesian optimization case, future work could explore this in the safety-critical case.

## 7 Proofs

In this section, we provide the proofs for Theorem 4.1 and Lemma 4.2.

Since we consider the surrogate function  $h(a, i)$  in (4), we obtain  $q + 1$  individual measurements at each iteration, each with individual noise. A measure of how difficult it

is to learn an accurate GP model of a function is given by the information capacity. This corresponds to the maximum amount of mutual information between a scalar function  $h$  and measurements  $\hat{\mathbf{h}}_{\mathcal{D}}$  at a set of parameters  $\mathcal{D}$  of size  $n$ . The measurements  $\hat{\mathbf{h}}_{\mathcal{D}}$  are corrupted by zero-mean, Gaussian noise. The information capacity is then defined as

$$\gamma_n := \max_{\mathcal{D} \subset \mathcal{A} \times \mathcal{I}, |\mathcal{D}|=n} \mathbf{I}(\hat{\mathbf{h}}_{\mathcal{D}}; h), \quad (38)$$

which is the maximum amount of information we can obtain about the function  $h$  from  $n$  measurements. The information gain is known to be sublinear in  $n$  for many commonly used kernels (Srinivas et al., 2012). Intuitively, the first samples for the GP model provide a lot of information, since each sample improves the prior significantly. After some iterations the domain  $\mathcal{A}$  is covered with samples in  $\mathcal{D}$ , so that additional samples are more correlated with previous data points in  $\mathcal{D}$ , rendering the samples less informative. The more prior information we encode in the GP prior, the less information can be gained from the same number of samples.

In our setting, we obtain  $|\mathcal{I}| = q + 1$  measurements at every iteration step  $n$ , each with different, independent noise. The mutual information with regards to these multiple measurements at parameters  $\bar{\mathcal{A}} \subset \mathcal{A}$  can be bounded with

$$\mathbf{I}(\hat{\mathbf{h}}_{\bar{\mathcal{A}} \times \mathcal{I}}; h) \leq \max_{\bar{\mathcal{A}} \subset \mathcal{A}, |\bar{\mathcal{A}}| \leq n} \mathbf{I}(\hat{\mathbf{h}}_{\bar{\mathcal{A}} \times \mathcal{I}}; h), \quad (39)$$

$$\leq \max_{\mathcal{D} \subset \mathcal{A} \times \mathcal{I}, |\mathcal{D}| \leq n|\mathcal{I}|} \mathbf{I}(\hat{\mathbf{h}}_{\mathcal{D}}; h), \quad (40)$$

$$= \gamma_{|\mathcal{I}|n}, \quad (41)$$

where  $\bar{\mathcal{A}} \times \mathcal{I}$  is the cartesian product that means we obtain one measurement for every function indexed by  $i \in \mathcal{I}$  at each parameter in  $\bar{\mathcal{A}}$ . The first inequality bounds the mutual information gained by Algorithm 1 by the worst-case mutual information, while the second inequality bounds this again by the worst-case mutual information when optimizing over the  $|\mathcal{I}|$  measurements at each iteration step individually. Intuitively, obtaining multiple optimal samples does not fundamentally change the properties of the information gain, but accelerates the rate at which information can be obtained in the worst case by  $|\mathcal{I}|$ .

In the following, we assume that  $h(\mathbf{a}, i)$  has bounded RKHS norm. Lemma 4.1 provides requirements for  $\beta_n$ , which will be used in the following to prove the results.

**Lemma 4.1** (based on Chowdhury and Gopalan (2017)) *Assume that  $h(\mathbf{a}, i)$  has RKHS norm bounded by  $B$  and that measurements are corrupted by  $\sigma$ -sub-Gaussian noise. If  $\beta_n^{1/2} = B + 4\sigma \sqrt{\gamma_{(n-1)|\mathcal{I}|} + 1 + \ln(1/\delta)}$ , then the following holds for all parameters  $\mathbf{a} \in \mathcal{A}$ , function indices  $i \in \mathcal{I}$ , and iterations  $n \geq 1$  jointly with probability at least  $1 - \delta$ :*

$$|h(\mathbf{a}, i) - \mu_{n-1}(\mathbf{a}, i)| \leq \beta_n^{1/2} \sigma_{n-1}(\mathbf{a}, i). \quad (18)$$

**Proof** Directly follows from Chowdhury and Gopalan (2017). The only difference is that we obtain  $|\mathcal{I}|$  measurements at every iteration, which causes the information capacity  $\gamma$  to grow at a faster rate.  $\square$

*Note* Where needed in the following lemmas, we implicitly assume that the assumptions of Lemma 4.1 hold, and that  $\beta_n$  is defined as above. Moreover, we use the definition for  $C(\mathbf{a}, i)$  as the intersection of the confidence intervals  $Q$  as in Sect. 4.1.

**Corollary 7.1** *For  $\beta_n$  as above, the following holds with probability at least  $1 - \delta$ :*

$$\forall n \geq 1, \forall i \in \mathcal{I}, \forall \mathbf{a} \in \mathcal{A}, h(\mathbf{a}, i) \in C_n(\mathbf{a}, i).$$

**Proof** From Lemma 4.1 we know that the true functions are contained in  $Q_n(\mathbf{a}, i)$  for all iterations  $n$  with probability at least  $1 - \delta$ . As a consequence, the true functions will be contained in the intersection of these sets with the same probability.  $\square$

Corollary 7.1 gives a choice of  $\beta_n$ , which ensures that all the function values of  $h$  are contained within their respective confidence intervals with high probability. In the remainder of the paper, we follow the outline of the proofs in Sui et al. (2015), but extended them to account for multiple constraints.

We start by showing the dynamics of important sets and functions. Most importantly, the upper confidence bounds are decreasing, lower confidence bounds increasing with the number of iterations, since the sets  $C_{n+1} \subseteq C_n$  for all iterations  $n$ .

**Lemma 7.1** *The following hold for any  $n \geq 1$ :*

- (i)  $\forall \mathbf{a} \in \mathcal{A}, \forall i \in \mathcal{I}, u_{n+1}^i(\mathbf{a}) \leq u_n^i(\mathbf{a}),$
- (ii)  $\forall \mathbf{a} \in \mathcal{A}, \forall i \in \mathcal{I}, l_{n+1}^i(\mathbf{a}) \geq l_n^i(\mathbf{a}),$
- (iii)  $\forall \mathbf{a} \in \mathcal{A}, \forall i \in \mathcal{I}, w_{n+1}(\mathbf{a}, i) \leq w_n(\mathbf{a}, i),$
- (iv)  $S_{n+1} \supseteq S_n \supseteq S_0,$
- (v)  $S \subseteq R \Rightarrow R_\epsilon(S) \subseteq R_\epsilon(R),$
- (vi)  $S \subseteq R \Rightarrow \tilde{R}_\epsilon(S) \subseteq \tilde{R}_\epsilon(R).$

**Proof** (i), (ii), and (iii) follow directly from their definitions and the definition of  $C_n(\mathbf{a})$ .

- (iv) Proof by induction. Consider the initial safe set,  $S_0$ . By definition of  $C_0$  we have for all  $\mathbf{a} \in S_0$  and  $i \in \mathcal{I}$  that

$$l_1^i(\mathbf{a}) - L\|\mathbf{a} - \mathbf{a}\| = l_1^i(\mathbf{a}) \geq l_0^i(\mathbf{a}) \geq 0.$$

It then follows from the definition of  $S_n$  that  $\mathbf{a} \in S_1$ .

For the induction step, assume that for some  $n \geq 2$ ,  $S_{n-1} \subseteq S_n$  and let  $\mathbf{a} \in S_n$ . This means that for all  $i \in \mathcal{I}_g$ ,  $\exists \mathbf{z}_i \in S_{n-1}$ ,  $l_n^i(\mathbf{z}_i) - L\|\mathbf{z}_i - \mathbf{a}\| \geq 0$  by the definition of the safe set. But, since  $S_{n-1} \subseteq S_n$ , this implies that  $\mathbf{z}_i \in S_n$ ,  $\forall i \in \mathcal{I}_g$ . Furthermore, by part (ii),  $l_{n+1}^i(\mathbf{z}_i) \geq l_n^i(\mathbf{z}_i)$ . Therefore, we conclude that for all  $i \in \mathcal{I}_g$ ,  $l_{n+1}^i(\mathbf{z}_i) - L\|\mathbf{z}_i - \mathbf{a}\| \geq 0$ , which implies that  $\mathbf{a} \in S_{n+1}$ .

- (v) Let  $\mathbf{a} \in R_\epsilon(S)$ . Then, by definition, for all  $i \in \mathcal{I}_g$ ,  $\exists \mathbf{z}_i \in S$ ,  $g_i(\mathbf{z}_i) - L\|\mathbf{z}_i - \mathbf{a}\| \geq 0$ . But, since  $S \subseteq R$ , it means that  $\mathbf{z}_i \in R \forall i \in \mathcal{I}_g$ , and, therefore,  $g_i(\mathbf{z}_i) - L\|\mathbf{z}_i - \mathbf{a}\| \geq 0$  for all  $i \in \mathcal{I}_g$  also implies that  $\mathbf{a} \in R_\epsilon(R)$ .
- (vi) This follows directly by repeatedly applying the result of part (v).

$\square$

Using the previous results, we start by showing that, after a finite number of iterations, the safe set has to expand if possible. As a first step, note that the set of expanders and maximizers are contained in each other as well if the safe set does not increase:

**Lemma 7.2** *For any  $n_1 \geq n_0 \geq 1$ , if  $S_{n_1} = S_{n_0}$ , then, for any  $n$ , such that  $n_0 \leq n < n_1$ , it holds that*

$$G_{n+1} \cup M_{n+1} \subseteq G_n \cup M_n.$$

**Proof** Given the assumption that  $S_n$  does not change, both  $G_{n+1} \subseteq G_n$  and  $M_{n+1} \subseteq M_n$  follow directly from the definitions of  $G_n$  and  $M_n$ . In particular, for  $G_n$ , note that for any  $\mathbf{a} \in S_n$ ,  $e_n^i(\mathbf{a})$  is decreasing in  $n$  for all  $i \in \mathcal{I}_g$ , since  $u_n^i(\mathbf{a})$  are decreasing in  $n$ . For  $M_n$ , note that  $\max_{\mathbf{a}' \in S_n} l_n^j(\mathbf{a}')$  is increasing in  $n$ , while  $u_n^i(\mathbf{a})$  is decreasing in  $n$  (see Lemma 7.1 (i), (ii)).  $\square$

When running the SAFE-Opt-MC algorithm, we repeatedly choose the most uncertain element from  $G_n$  and  $M_n$ . Since these sets are contained in each other if the safe set does not expand, we gain more information about these sets with each sample. Since the information gain is bounded, this allows us to bound the uncertainty in terms of the information gain over the entire set:

**Lemma 7.3** *For any  $n_1 \geq n_0 \geq 1$ , if  $S_{n_1} = S_{n_0}$  and  $\tilde{c} := 8/\log(1 + \sigma^{-2})$ , then, for any  $n$ , such that  $n_0 \leq t \leq n_1$ , it holds for all  $i \in \mathcal{I}$  that*

$$w_n(\mathbf{a}_n, i) \leq \sqrt{\frac{\tilde{c} \beta_n \gamma_{|\mathcal{I}|n}}{n - n_0}}.$$

**Proof** Given Lemma 7.2, the definition of  $\mathbf{a}_n := \operatorname{argmax}_{\mathbf{a} \in G_n \cup M_n} (w_n(\mathbf{a}))$ , and the fact that,  $w_n^i(\mathbf{a}_n) \leq 2\beta_n^{1/2} \max_{i \in \mathcal{I}} \sigma_{n-1}(\mathbf{a}_n, i) = 2\beta_n^{1/2}(\mathbf{a}_n, i_n)$ , the proof is completely analogous to that of Lemma 5.3 by Srinivas et al. (2012). We only highlight the main differences here, which results from having several functions.

$$w_n^i(\mathbf{a}_n) \leq 2\beta_n^{1/2} \max_{i \in \mathcal{I}} \sigma_{n-1}(\mathbf{a}_n, i), \quad (42)$$

which following (Srinivas et al., 2012, Lemma 5.4) leads to

$$\sum_{j=1}^n w_j^2(\mathbf{a}_j, i_j) \leq \beta_{|\mathcal{I}|n}^{1/2} \mathbf{I}(\hat{\mathbf{h}}_{\bar{\mathcal{D}}_n}; h), \quad (43)$$

where  $\bar{\mathcal{D}}_n = \{\mathbf{a}_n, i_n\}$ . Now using monotonicity of the mutual information, we have that

$$\sum_{j=1}^n w_j^2(\mathbf{a}_j, i_j) \leq \tilde{c} \beta_{|\mathcal{I}|n}^{1/2} \mathbf{I}(\hat{\mathbf{h}}_{\mathcal{D}_n \times \mathcal{I}}; h), \quad (44)$$

$$\leq \tilde{c} \beta_{|\mathcal{I}|n}^{1/2} \gamma_{|\mathcal{I}|n} \quad (45)$$

by (41).  $\square$

**Corollary 7.2** For any  $n \geq 1$ , if  $\tilde{c}$  is defined as above,  $N_n$  is the smallest positive integer satisfying  $\frac{N_n}{\beta_{n+N_n} \gamma_{|\mathcal{I}|(n+N_n)}} \geq \frac{\tilde{c}}{\epsilon^2}$ , and  $S_{n+N_n} = S_n$ , then, for any  $\mathbf{a} \in G_{n+N_n} \cup M_{n+N_n}$ , and for all  $i \in \mathcal{I}$  it holds that

$$w_{n+N_n}(\mathbf{a}, i) \leq \epsilon.$$

*Note* Where needed in the following lemmas, we assume that  $\tilde{c}$  and  $N_n$  are defined as above.

That is, after a finite number of evaluations  $N_n$  the most uncertain element within these sets is at most  $\epsilon$ . Given that the reachability operator in (8) is defined in terms of the same accuracy, it allows us to show that after at most  $N_n$  evaluations, the safe set has to increase unless it is impossible to do so:

**Lemma 7.4** For any  $n \geq 1$ , if  $\bar{R}_\epsilon(S_0) \setminus S_n \neq \emptyset$ , then  $R_\epsilon(S_n) \setminus S_n \neq \emptyset$ .

**Proof** Assume, to the contrary, that  $R_\epsilon(S_n) \setminus S_n = \emptyset$ . By definition,  $R_\epsilon(S_n) \supseteq S_n$ , therefore  $R_\epsilon(S_n) = S_n$ . Iteratively applying  $R_\epsilon$  to both sides, we get in the limit  $\bar{R}_\epsilon(S_n) = S_n$ . But then, by Lemma 7.1 (iv) and (vi), we get

$$\bar{R}_\epsilon(S_0) \subseteq \bar{R}_\epsilon(S_n) = S_n, \quad (46)$$

which contradicts the lemma's assumption that  $\bar{R}_\epsilon(S_0) \setminus S_n \neq \emptyset$ .  $\square$

**Lemma 7.5** For any  $n \geq 1$ , if  $\bar{R}_\epsilon(S_0) \setminus S_n \neq \emptyset$ , then the following holds with probability at least  $1 - \delta$ :

$$S_{n+N_n} \supsetneq S_n.$$

**Proof** By Lemma 7.4, we get that,  $R_\epsilon(S_n) \setminus S_n \neq \emptyset$ . Equivalently, by definition, for all  $i \in \mathcal{I}_g$

$$\exists \mathbf{a} \in R_\epsilon(S_n) \setminus S_n, \exists \mathbf{z}_i \in S_n : g_i(\mathbf{z}_i) - \epsilon - L\|\mathbf{z}_i - \mathbf{a}\| \geq 0. \quad (47)$$

Now, assume, to the contrary, that  $S_{n+N_n} = S_n$  (see Lemma 7.1 (iv)), which implies that  $\mathbf{a} \in \mathcal{A} \setminus S_{n+N_n}$  and  $\mathbf{z}_i \in S_{n+N_n} \forall i \in \mathcal{I}_g$ . Then, we have for all  $i \in \mathcal{I}_g$

$$\begin{aligned} u_{n+N_n}^i(\mathbf{z}_i) - L\|\mathbf{z}_i - \mathbf{a}\| &\geq g_i(\mathbf{z}_i) - L\|\mathbf{z} - \mathbf{a}\| && \text{by Lemma 4.1} \\ &\geq g_i(\mathbf{z}_i) - \epsilon - L\|\mathbf{z} - \mathbf{a}\| \\ &\geq 0. && \text{by (47)} \end{aligned}$$

Therefore, by definition,  $e_{n+N_n}(\mathbf{z}_i) > 0$ , which implies  $\mathbf{z}_i \in G_{n+N_n}, \forall i \in \mathcal{I}_g$ .

Finally, since  $S_{n+N_n} = S_n$  and  $\mathbf{z}_i \in G_{n+N_n} \forall i \in \mathcal{I}_g$ , we know that for all  $i \in \mathcal{I}$ ,  $w_{n+N_n}(\mathbf{a}', i) \leq \epsilon$ . (Corollary 7.2). Hence, for all  $i \in \mathcal{I}_g$ ,

$$\begin{aligned} l_{n+N_n}^i(\mathbf{z}_i) - L\|\mathbf{z}_i - \mathbf{a}\| &\geq g_i(\mathbf{z}_i) - w(\mathbf{z}_i, i) - L\|\mathbf{a} - \mathbf{z}_i\| && \text{by Lemma 4.1} \\ &\geq g_i(\mathbf{z}) - \epsilon - L\|\mathbf{a} - \mathbf{z}_i\| && \text{by Corollary 7.2} \\ &\geq 0 && \text{by (47).} \end{aligned}$$

This means we get  $\mathbf{a} \in S_{n+N_n}$ , which is a contradiction.  $\square$

Intuitively, repeatedly applying the previous result leads to full safe exploration within a finite domain  $\mathcal{A}$ . In particular, it follows that if  $S_{n+N_n} = S_n$ , then the safely reachable set has been fully explored to the desired accuracy. From this it follows, that the pessimistic estimate in (17) is also  $\epsilon$ -close to the optimum value within the safely reachable set,  $\bar{R}_\epsilon(S_0)$ :

**Lemma 7.6** *For any  $n \geq 1$ , if  $S_{n+N_n} = S_n$ , then the following holds with probability at least  $1 - \delta$ :*

$$f(\mathbf{a}_{n+N_n}) \geq \max_{\mathbf{a} \in \bar{R}_\epsilon(S_0)} f(\mathbf{a}) - \epsilon.$$

**Proof** Let  $\mathbf{a}^* := \operatorname{argmax}_{\mathbf{a} \in S_{n+N_n}} f(\mathbf{a})$ . Note that  $\mathbf{a}^* \in M_{n+N_n}$ , since

$$\begin{aligned} u_{n+N_n}^f(\mathbf{a}^*) &\geq f(\mathbf{a}^*) && \text{by Lemma 4.1} \\ &\geq f(\mathbf{a}) && \text{by definition of } \mathbf{a}^* \\ &\geq l_{n+N_n}^f(\mathbf{a}) && \text{by Lemma 4.1} \\ &\geq \max_{\mathbf{a} \in S_{n+N_n}} l_{n+N_n}^f(\mathbf{a}). && \text{by definition of } \mathbf{a} \end{aligned}$$

We will first show that  $f(\mathbf{a}_{n+N_n}) \geq f(\mathbf{a}^*) - \epsilon$ . Assume, to the contrary, that

$$f(\mathbf{a}_{n+N_n}) < f(\mathbf{a}^*) - \epsilon. \quad (48)$$

Then, we have

$$\begin{aligned} l_{n+N_n}^f(\mathbf{a}^*) &\leq l_{n+N_n}^f(\mathbf{a}) && \text{by definition of } \mathbf{a} \\ &\leq f(\mathbf{a}) && \text{by Lemma 4.1} \\ &< f(\mathbf{a}^*) - \epsilon && \text{by (48)} \\ &\leq u_{n+N_n}^f(\mathbf{a}^*) - \epsilon && \text{by Lemma 4.1} \\ &\leq l_{n+N_n}^f(\mathbf{a}^*), && \text{by Corollary 7.2 and } \mathbf{a}^* \in M_{n+N_n} \end{aligned}$$

which is a contradiction.

Finally, since  $S_{n+N_n} = S_n$ , Lemma 7.5 implies that  $\bar{R}_\epsilon(S_0) \subseteq S_n = S_{n+N_n}$ . Therefore,

$$\begin{aligned} \max_{\mathbf{a} \in \bar{R}_\epsilon(S_0)} f(\mathbf{a}) - \epsilon &\leq \max_{\mathbf{a} \in S_{n+N_n}} f(\mathbf{a}) - \epsilon && \bar{R}_\epsilon(S_0) \subseteq S_{n+N_n} \\ &= f(\mathbf{a}^*) - \epsilon && \text{by definition of } \mathbf{a}^* \\ &\leq f(\mathbf{a}_{n+N_n}) && \text{proven above.} \end{aligned}$$

□

**Corollary 7.3** *For any  $n \geq 1$ , if  $S_{n+N_n} = S_n$ , then the following holds with probability at least  $1 - \delta$ :*

$$\forall n' \geq 0, f(\mathbf{a}_{n+N_n+n'}) \geq \max_{\mathbf{a} \in \bar{R}_\epsilon(S_0)} f(\mathbf{a}) - \epsilon.$$

**Proof** This is a direct consequence of the proof of the preceding lemma, combined with the facts that both  $S_{n+N_n+n'}$  and  $l_{n+N_n+n'}^f(\mathbf{a}_{n+N_n+n'})$  are increasing in  $n'$  (by Lemma 7.1 (iv) and (ii) respectively), which imply that  $\max_{\mathbf{a} \in S_{n+N_n+n'}} l_{n+N_n+n'}^f(\mathbf{a})$  can only increase in  $n'$ . □

Moreover, since we know the true function is contained within the confidence intervals, we cannot go beyond the safe set if we knew the function perfectly everywhere,  $\bar{R}_0$ :

**Lemma 7.7** *For any  $n \geq 0$ , the following holds with probability at least  $1 - \delta$ :*

$$S_n \subseteq \bar{R}_0(S_0).$$

**Proof** Proof by induction. For the base case,  $n = 0$ , we have by definition that  $S_0 \subseteq \bar{R}_0(S_0)$ .

For the induction step, assume that for some  $n \geq 1$ ,  $S_{n-1} \subseteq \bar{R}_0(S_0)$ . Let  $\mathbf{a} \in S_n$ , which, by definition, means that for all  $i \in \mathcal{I}_g \exists \mathbf{z}_i \in S_{n-1}$ , such that

$$\begin{aligned} l_n^i(\mathbf{z}_i) - L\|\mathbf{z}_i - \mathbf{a}\| &\geq 0 \\ \Rightarrow g_i(\mathbf{z}_i) - L\|\mathbf{z}_i - \mathbf{a}\| &\geq 0. \quad \text{by Lemma 4.1} \end{aligned}$$

Then, by definition of  $\bar{R}_0$  and the fact that  $\mathbf{z}_i \in \bar{R}_0(S_0)$  for all  $i \in \mathcal{I}_g$ , it follows that  $\mathbf{a} \in \bar{R}_0(S_0)$ .  $\square$

The previous results is enough to show that we eventually explore the full safe set by repeatedly applying Lemma 7.5:

**Lemma 7.8** *Let  $n^*$  be the smallest integer, such that  $n^* \geq |\bar{R}_0(S_0)|T_{n^*}$ . Then, there exists  $n_0 \leq n^*$ , such that  $S_{n_0+T_{n_0}} = S_{n_0}$ .*

**Proof** Assume, to the contrary, that for any  $n \leq n^*$ ,  $S_n \subsetneq S_{n+T_n}$ . (By Lemma 7.1 (iv), we know that  $S_n \subseteq S_{n+T_n}$ .) Since  $N_n$  is increasing in  $n$ , we have

$$S_0 \subsetneq S_{n_0} \subseteq S_{T_{n^*}} \subsetneq S_{T_{n^*}+T_{T_{n^*}}} \subseteq S_{2T_{n^*}} \subsetneq \dots,$$

which implies that, for any  $0 \leq k \leq |\bar{R}_0(S_0)|$ , it holds that  $|S_{kT_{n^*}}| > k$ . In particular, for  $k^* := |\bar{R}_0(S_0)|$ , we get

$$|S_{k^*T}| > |\bar{R}_0(S_0)|$$

which contradicts  $S_{k^*T} \subseteq \bar{R}_0(S_0)$  by Lemma 7.7.  $\square$

**Corollary 7.4** *Let  $n^*$  be the smallest integer, such that  $\frac{n^*}{\beta_{n^*}\gamma_{|\mathcal{I}|n^*}} \geq \frac{\tilde{c}|\bar{R}_0(S_0)|}{\epsilon^2}$ . Then, there exists  $n_0 \leq n^*$ , such that  $S_{n_0+T_{n_0}} = S_{n_0}$ .*

**Proof** This is a direct consequence of combining Lemma 7.8 and Corollary 7.2.  $\square$

Since we showed that we completely explore the safe set and that we remain safe throughout the exploration procedure, we are ready to state the main results:

**Lemma 7.9** *If  $h$  is  $L$ -Lipschitz continuous, then, for any  $n \geq 0$ , the following holds with probability at least  $1 - \delta$  for all  $i \in \mathcal{I}_g$ :*

$$\forall \mathbf{a} \in S_n, g_i(\mathbf{a}) \geq 0.$$

**Proof** We will prove this by induction. For the base case  $n = 0$ , by definition, for any  $\mathbf{a} \in S_0$  and  $i \in \mathcal{I}_g$ ,  $g_i(\mathbf{a}) \geq 0$ .



For the induction step, assume that for some  $n \geq 1$ , for any  $\mathbf{a} \in S_{n-1}$  and for all  $i \in \mathcal{I}_g$ ,  $g_i(\mathbf{a}) \geq 0$ . Then, for any  $\mathbf{a} \in S_n$ , by definition, for all  $i \in \mathcal{I}_g$ ,  $\exists \mathbf{z}_i \in S_{n-1}$ ,

$$\begin{aligned} 0 &\leq l_n^i(\mathbf{z}_i) - L\|\mathbf{z}_i - \mathbf{a}\| \\ &\leq g_i(\mathbf{z}_i) - L\|\mathbf{z}_i - \mathbf{a}\| \quad \text{by Lemma 4.1} \\ &\leq g_i(\mathbf{a}) \quad \text{by } L\text{-Lipschitz-continuity.} \end{aligned}$$

□

**Theorem 4.1** Assume that  $h(\mathbf{a}, i)$  has bounded norm in an RKHS and that the measurement noise is  $\sigma$ -sub-Gaussian. Also, assume that  $S_0 \neq \emptyset$  and  $g_i(\mathbf{a}) \geq 0$  for all  $\mathbf{a} \in S_0$  and  $i \in \mathcal{I}_g$ . Choose  $\beta_n$  as in Lemma 4.1, define  $\hat{\mathbf{a}}_n$  as in (17), and let  $n^*(\epsilon, \delta)$  be the smallest positive integer satisfying

$$\frac{n^*}{\beta_{n^*} \gamma_{|\mathcal{I}|n^*}} \geq \frac{\tilde{c}(|\bar{R}_0(S_0)| + 1)}{\epsilon^2}, \quad (20)$$

where  $\tilde{c} = 8/\log(1 + \sigma^{-2})$ . For any  $\epsilon > 0$  and  $\delta \in (0, 1)$ , when running Algorithm 1 the following inequalities jointly hold with probability at least  $1 - \delta$ :

1. *Safety*:  $\forall n \geq 1, \forall i \in \mathcal{I}_g : g_i(\mathbf{a}_n) \geq 0$
2. *Optimality*:  $\forall n \geq n^*, f(\hat{\mathbf{a}}_n) \geq f_\epsilon^* - \epsilon$

**Proof** The first part of the theorem is a direct consequence of Lemma 7.9. The second part follows from combining Corollary 7.3 and Corollary 7.4. □

**Acknowledgements** This research was supported in part by SNSF Grant 200020\_159557, NSERC Grant RGPIN-2014-04634, and the Connaught New Researcher Award.

**Funding** Open Access funding provided by ETH Zurich.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Achiam, J., Held, D., Tamar, A., & Abbeel, P. (2017). Constrained policy optimization. In *Proceedings of the international conference on machine learning (ICML)*, 1705.10528.
- Akametalu, A. K., Kaynama, S., Fisac, J. F., Zeilinger, M. N., Gillula, J. H., & Tomlin, C. J. (2014). Reachability-based safe learning with Gaussian processes. In *Proceedings of the IEEE conference on decision and control (CDC)* (pp. 1424–1431).
- Álvarez, M. A., Rosasco, L., & Lawrence, N. D. (2012). Kernels for vector-valued functions: A review. *Foundations and Trends in Machine Learning*, 4(3), 195–266.

- Aswani, A., Gonzalez, H., Sastry, S. S., & Tomlin, C. (2013). Provably safe and robust learning-based model predictive control. *Automatica*, 49(5), 1216–1226.
- Berkenkamp, F., Krause, A., Angela, P., & Schoellig (2016a). Bayesian optimization with safety constraints: Safe and automatic parameter tuning in robotics. [arXiv:160204450](https://arxiv.org/abs/160204450) [csRO]
- Berkenkamp, F., & Schoellig, A. P. (2015). Safe and robust learning control with Gaussian processes. In *Proceedings of the European control conference (ECC)* (pp. 2501–2506).
- Berkenkamp, F., Schoellig, A. P., & Krause, A. (2016b). Safe controller optimization for quadrotors with Gaussian processes. In *IEEE international conference on robotics and automation (ICRA)* (pp. 493–496).
- Berkenkamp, F., Turchetta, M., Schoellig, A. P., & Krause, A. (2017). Safe model-based reinforcement learning with stability guarantees. In *Neural information processing systems (NeurIPS)* (pp. 908–918).
- Boyd, S., Boyd, S. P., & Vandenberghe, L. (2004). *Convex optimization*. Cambridge University Press.
- Bull, A. D. (2011). Convergence rates of efficient global optimization algorithms. *Journal of Machine Learning Research*, 12, 2879–2904.
- Calandra, R., Gopalan, N., Seyfarth, A., Peters, J., & Deisenroth, M. P. (2014a). Bayesian gait optimization for bipedal locomotion. In *Learning and intelligent optimization* (pp 274–290). Springer.
- Calandra, R., Seyfarth, A., Peters, J., Deisenroth, M. P. (2014b). An experimental comparison of Bayesian optimization for bipedal locomotion. In *2014 IEEE international conference on robotics and automation (ICRA)* (pp. 1951–1958).
- Chowdhury, S. R., & Gopalan, A. (2017). On kernelized multi-armed bandits. In *Proceedings of the 34th international conference on machine learning, PMLR, proceedings of machine learning research* (Vol. 70, pp. 844–853).
- Christmann, A., & Steinwart, I. (2008). *Support vector machines. Information science and statistics*. Springer.
- Dalal, G., Dvijotham, K., Vecerik, M., Hester, T., Paduraru, C., & Tassa, Y. (2018). Safe exploration in continuous action spaces. [arXiv preprint arXiv:180108757](https://arxiv.org/abs/180108757)
- Djolonga, J., Krause, A., & Cevher, V. (2013). High-dimensional Gaussian process bandits. *Advances in Neural Information Processing Systems*, 26, 1025–1033.
- Duivenvoorden, R. R., Berkenkamp, F., Carion, N., Krause, A., & Schoellig, A. P. (2017). Constrained Bayesian optimization with particle swarms for adaptive controller tuning. In *Proceedings of the IFAC (international federation of automatic control) world congress* (pp. 12306–12313).
- Duvenaud, D. K., Nickisch, H., & Rasmussen, C. E. (2011). Additive Gaussian processes. *Advances in Neural Information Processing Systems*, 24, 226–234.
- Eldar, Y. C., & Kutyniok, G. (2012). *Compressed sensing: Theory and applications*. Cambridge University Press.
- Garcia, J., & Fernández, F. (2012). Safe exploration of state and action spaces in reinforcement learning. *Journal of Artificial Intelligence Research*, 45, 515–564.
- Gelbart, M.A., Snoek, J., Adams, R.P. (2014). Bayesian optimization with unknown constraints. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI)* (pp. 250–259).
- Ghosal, S., & Roy, A. (2006). Posterior consistency of Gaussian process prior for nonparametric binary regression. *The Annals of Statistics*, 34(5), 2413–2429.
- Hans, A., Schneegaß, D., Schäfer, A. M., & Udluft, S. (2008). Safe exploration for reinforcement learning. In *ESANN* (pp. 143–148). Citeseer.
- Jones, D. R. (2001). A taxonomy of global optimization methods based on response surfaces. *Journal of Global Optimization*, 21(4), 345–383.
- Kober, J., & Peters, J. (2013). Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11), 1238–1274.
- Krause, A., & Ong, C. S. (2011). Contextual Gaussian process bandit optimization. In *Proceedings of neural information processing systems (NIPS)* (pp. 2447–2455).
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., & Wierstra, D. (2015). Continuous control with deep reinforcement learning. [arXiv:150902971](https://arxiv.org/abs/150902971) [cs, stat] 1509.02971
- Lizotte, D. J., Wang, T., Bowling, M. H., & Schuurmans, D. (2007). Automatic gait optimization with Gaussian process regression. In *Proceedings of the twentieth international joint conference on artificial intelligence (IJCAI)* (Vol. 7, pp. 944–949).
- Lupashin, S., Hehn, M., Mueller, M. W., Schoellig, A. P., Sherback, M., & D’Andrea, R. (2014). A platform for aerial robotics research and demonstration: The Flying Machine Arena. *Mechatronics*, 24(1), 41–54.
- Marco, A., Berkenkamp, F., Hennig, P., Schoellig, A. P., Krause, A., Schaal, S., & Trimpe, S. (2017). Virtual vs. real: Trading off simulations and physical experiments in reinforcement learning with Bayesian

- optimization. In *Proceedings of the IEEE international conference on robotics and automation (ICRA)* (pp. 1557–1563).
- Mockus, J. (2012). *Bayesian approach to global optimization: Theory and applications*. Springer
- Moldovan, T. M., & Abbeel, P. (2012). Safe exploration in Markov decision processes. In *Proceedings of the international conference on machine learning (ICML)* (pp. 1711–1718).
- Ostafew, C. J., Schoellig, A. P., & Barfoot, T. D. (2016). Robust constrained learning-based NMPC enabling reliable mobile robot path tracking. *The International Journal of Robotics Research (IJRR)*, 35(13), 1547–1536.
- Peters, J., & Schaal, S. (2006). Policy gradient methods for robotics. In *Proceedings of the IEEE/RSJ international conference on intelligent robots and systems* (pp. 2219–2225).
- Peters, J., & Schaal, S. (2008). Reinforcement learning of motor skills with policy gradients. *Neural Networks*, 21(4), 682–697.
- Rasmussen, C. E., & Williams, C. K. (2006). *Gaussian processes for machine learning*. Cambridge MA: MIT Press.
- Schaal, S., & Atkeson, C. G. (2010). Learning control in robotics. *IEEE Robotics & Automation Magazine*, 17(2), 20–29.
- Schoellig, A., Wiltscche, C., & D'Andrea, R. (2012). Feed-forward parameter identification for precise periodic quadcopter motions. In *Proceedings of the American control conference (ACC)* (pp. 4313–4318).
- Schoellig, A. P., Hehn, M., Lupashin, S., & D'Andrea, R. (2011). Feasibility of motion primitives for choreographed quadcopter flight. In *Proceedings of the American control conference (ACC)* (pp. 3843–3849).
- Schreiter, J., Nguyen-Tuong, D., Eberts, M., Bischoff, B., Markert, H., & Toussaint, M. (2015). Safe exploration for active learning with Gaussian processes. In *Machine learning and knowledge discovery in databases* (Vol. 9286, pp. 133–149). Springer International Publishing.
- Solak, E., Murray-Smith, R., Leithead, W. E., Leith, D. J., & Rasmussen, C. E. (2003). Derivative observations in Gaussian process models of dynamic systems. In S. Becker, S. Thrun, & K. Obermayer (eds) *Proceedings of neural information processing systems (NIPS)* (pp. 1057–1064). MIT Press.
- Srinivas, N., Krause, A., Kakade, S. M., & Seeger, M. (2012). Gaussian process optimization in the bandit setting: No regret and experimental design. *IEEE Transactions on Information Theory*, 58(5), 3250–3265 ((0912.3995)).
- Sui, Y., Gotovos, A., Burdick, J. W., & Krause, A. (2015). Safe exploration for optimization with Gaussian processes. In *Proceedings of the international conference on machine learning (ICML)* (pp. 997–1005).
- Sui, Y., Zhuang, B. J., & Yue, Y. (2018). Stagewise safe Bayesian optimization with Gaussian processes. In *International conference on machine learning* (pp. 4781–4789).
- Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning: An introduction*. MIT press.
- Tesch, M., Schneider, J., & Choset, H. (2011). Using response surfaces and expected improvement to optimize snake robot gait parameters. In *Proceedings of the IEEE/RSJ international conference on intelligent robots and systems (IROS)* (pp. 1069–1074).
- The GPy authors (2012) GPy: A Gaussian process framework in Python. <https://github.com/SheffieldML/GPy>
- Turchetta, M., Berkenkamp, F., & Krause, A. (2016). Safe exploration in finite markov decision processes with Gaussian processes. In *Neural information processing systems (NeurIPS)* (pp. 4305–4313).
- Turchetta, M., Berkenkamp, F., & Krause, A. (2019). Safe exploration for interactive machine learning. In *Neural information processing systems (NeurIPS)*.
- Wang, Z., Zoghi, M., Hutter, F., Matheson, D., & De Freitas, N. (2013). Bayesian optimization in high dimensions via random embeddings. In *Proceedings of the international joint conference on artificial intelligence (IJCAI)* (pp. 1778–1784). AAAI Press.
- Zhou, K., & Doyle, J. C. (1998). *Essentials of robust control*, (Vol. 104). Prentice Hall.