

Computing Optimal Joint Chance Constrained Control Policies

Niklas Schmid, Marta Fochesato, Sarah H.Q. Li, Tobias Sutter, John Lygeros

Abstract— We consider the problem of optimally controlling stochastic, Markovian systems subject to joint chance constraints over a finite-time horizon. For such problems, standard Dynamic Programming is inapplicable due to the time correlation of the joint chance constraints, which calls for non-Markovian, and possibly stochastic, policies. Hence, despite the popularity of this problem, solution approaches capable of providing provably-optimal and easy-to-compute policies are still missing. We fill this gap by introducing an augmented binary state to the system dynamics, allowing us to characterize the optimal policies and propose a Dynamic Programming based solution method. Our analysis provides a deep insight into the impact of joint chance constraints on the optimal control policies.

Index Terms— Stochastic Optimal Control, Joint Chance Constrained Programming, Dynamic Programming.

I. INTRODUCTION

Many real-world control applications come with the requirement for safety certificates, for example in air-traffic control [1], [2], self-driving cars [3], robot path planning [4] or medicine [5]. In the wake of the rapid emergence of ever more complex safety-critical control problems involving stochastic systems, there is a growing need for optimal control tools outputting actions that are provably safe and easy to compute. Safety is commonly defined via a predefined set of safe states, in the sense that any state trajectory leaving the safe set during the control task is considered unsafe. For stochastic systems, safety can only be guaranteed up to some predefined probability. This leads to the definition of so-called joint chance constraints, i.e., constraints that bound the probability to be unsafe with respect to the entire state trajectory over a finite-time horizon. This is in contrast to stage-wise chance constrained formulations, which bound the probability to be unsafe at every individual time step. Such guarantees are especially popular for infinite-horizon problems as typically addressed in Model Predictive Control since the probability of remaining safe at all time-steps over the infinite trajectory is zero under mild assumptions on the system dynamics [6]–[9]. For the finite-time setting, however, as considered in this paper, we argue that it is generally more meaningful to define safety as a guarantee on the entire trajectory.

N. Schmid, M. Fochesato, S.H.Q. Li and J. Lygeros are with the Automatic Control Laboratory (IfA), ETH Zürich, 8092 Zürich, Switzerland {nschmid, mfochesato, huilih, jlygeros}@ethz.ch

T. Sutter is with the Department of Computer Science, University of Konstanz, 78457 Konstanz, Germany tobias.sutter@uni-konstanz.de

Work supported by the European Research Council under the Horizon 2020 Advanced under Grant 787845 (OCAL).

Despite the ubiquity of joint chance constrained problems [4], [6], [10]–[16], their practical deployment is hindered by their intractability due to (i) the difficulty in evaluating multivariate integrals to check the feasibility of a candidate solution [17]; (ii) the non-convex feasible region described by these constraints; and (iii) the time correlation introduced by the constraints. While (i) and (ii) equally affect stage-wise and joint chance constraints, (iii) is specific to joint chance constraints only. In particular, it is this time correlation which breaks the Markovian structure of the problem [6], forcing us to consider the full problem dimension at once.

To circumvent this problem, approximations of joint chance constraints have been proposed in the literature. Perhaps the best known is [4], that exploits Boole’s inequality to break the time correlation, leading to an inner approximation of the constraint set. More recently, [10] proposes to augment the state vector with a function space to fit the standard Dynamic Programming (DP) format. However, the resulting formulation is computationally challenging, even for one-dimensional examples, hindering the development of provably-convergent algorithms. Another recent approach relies on approximating the stochastic system dynamics using kernel distribution embeddings, allowing for data-based scenarios [15]. Motivated by portfolio management applications, [18] studies stochastic optimal control problems with a chance constraint only on the final state, circumventing the problem of time correlations among stages. On a broader scale, joint chance constrained problems have also been explored in combination with Model Predictive Control, approximating the problem using Boole’s inequality and similar bounds [12], [13], [19] or sampling [11], [14], [16]. Similar problems have further been addressed in the reinforcement learning community in the form of general constrained Markov Decision Processes [20]–[24].

To the best of the authors knowledge, no tractable solution to joint chance constrained optimal control problems has been proposed that is not reliant on conservative approximations. This paper aims to develop a deeper understanding of safety-constrained stochastic optimal control problems where safety is enforced via the use of joint chance constraints over the finite-time horizon of the mission. Specifically, our analysis establishes a novel DP recursion leading to a failure-aware policy that achieves an optimal trade-off between performance and safety. We highlight the following key results of this paper.

- (i) **Characterization of the optimal policies.** We propose a DP recursion to handle joint chance constraints in stochastic optimal control problems. The recursion is based on a suitable state augmentation that allows us

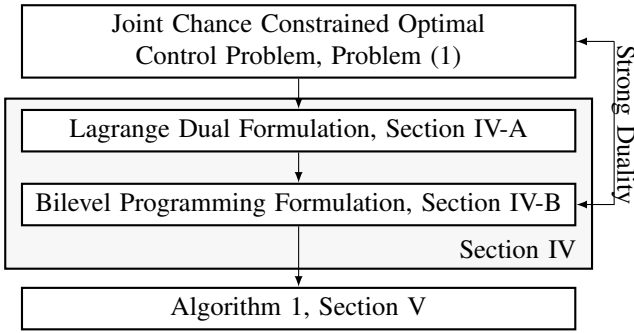


Fig. 1: Graphical representation of the paper structure.

to turn the original problem into an equivalent MDP. Based on the equivalence of stochastic causal policies and mixtures of deterministic Markov policies for the considered problem setting, we derive tractable solution methods based on the latter policy class.

- (ii) **Computation of the optimal policies.** We propose to solve joint chance constrained problems via their Lagrangian dual. We prove strong duality independent of the convexity of the cost functions and safe set, and exploit this fact to develop a bilevel optimization framework to solve the dual problem. Our approach yields an optimal trade-off between performance and safety.
- (iii) **Behaviour of joint chance constrained controllers.** We uncover and convey a clear intuition of the behaviour of joint chance constrained policies, which turns out to be “controversial” for many applications. The reason is that the associated optimal policies might advertise inputs that are likely to cause safety violations when every safe input yields high expected costs.

Unlike [4], [11]–[13] our method does not rely on conservative approximations, unlike [10] it is tractable and unlike [14]–[16], which rely on sampling techniques, it provides provably-optimal solutions.

The rest of the paper is organized as follows. In Section II we provide the problem formulation, in Section III we discuss the necessary background. In Section IV, we present our novel joint chance constrained Dynamic Programming scheme (see Fig. 1). In Section V we present our algorithmic solution and evaluate it numerically in Section VI. We conclude with a summary and outlook in Section VII.

Notation. We denote by $\mathbb{1}_{\mathcal{A}}(x)$ the indicator function of a set \mathcal{A} , where $\mathbb{1}_{\mathcal{A}}(x) = 1$ if $x \in \mathcal{A}$ and $\mathbb{1}_{\mathcal{A}}(x) = 0$ otherwise. Given two sets \mathcal{X}, \mathcal{Y} , the difference between them is denoted by $\mathcal{X} \setminus \mathcal{Y} = \{x \in \mathcal{X} : x \notin \mathcal{Y}\}$, while the complement of a set \mathcal{X} is denoted as \mathcal{X}^c . We denote by $[N]$ the set $\{0, 1, \dots, N\}$ and $\mathbb{R}_{\geq 0}$ the non-negative reals. Further, \wedge and \vee symbolize the logical conjunction and disjunction, respectively.

II. PROBLEM FORMULATION

We define a safety-constrained discrete-time stochastic system over a finite time-horizon as a tuple $(\mathcal{X}, \mathcal{U}, T, \ell_{0:N}, \mathcal{A})$, where the state space \mathcal{X} and the input space \mathcal{U} are Borel subsets of complete separable metric spaces equipped respectively with σ -algebras $\mathcal{B}(\mathcal{X})$ and $\mathcal{B}(\mathcal{U})$. Given a state

$x_k \in \mathcal{X}$ and an input $u_k \in \mathcal{U}$, the Borel-measurable stochastic kernel $T : \mathcal{B}(\mathcal{X}) \times \mathcal{X} \times \mathcal{U} \rightarrow [0, 1]$ describes the stochastic state evolution, leading to $x_{k+1} \sim T(\cdot | x_k, u_k)$. Additionally, $\ell_k : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}_{\geq 0}, k \in [N - 1]$ and $\ell_N : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ denote measurable, non-negative functions called stage and terminal cost, respectively, which are incurred at every time-step $k \in [N - 1]$ and at terminal time $N \in \mathbb{N}$. Finally, let the set $\mathcal{A} \subseteq \mathcal{B}(\mathcal{X})$ denote a safe set. We define a trajectory to be safe, if $x_{0:N} \in \mathcal{A}$. By abuse of notation, we interpret $x_{0:N} \in \mathcal{A}$ to mean $x_k \in \mathcal{A}$ for all $k \in [N]$.

For $k \in [N - 1]$ we define the space of histories up to time k recursively as $\mathcal{H}_k = \mathcal{X} \times \mathcal{U} \times \mathcal{H}_{k-1}$, with $\mathcal{H}_0 = \mathcal{X}$; a generic element $h_k \in \mathcal{H}_k$ is of the form $h_k = (x_0, u_0, x_1, u_1, \dots, x_{k-1}, u_{k-1}, x_k)$. We consider the following classes of policies:

A stochastic policy is a sequence $\pi = (\mu_0, \dots, \mu_{N-1})$ of Borel-measurable stochastic kernels $\mu_k, k \in [N - 1]$, that, given h_k , assigns a probability measure $\mu_k(\cdot, h_k)$ on the set $\mathcal{B}(\mathcal{U})$. A deterministic policy is the special case $\pi_{\text{d}} = (\mu_0, \dots, \mu_{N-1})$ where $\mu_k : \mathcal{H}_k \rightarrow \mathcal{U}, k \in [N - 1]$ are simply measurable maps. A policy is further called Markov if $\mu_k(\cdot | h_k) = \mu_k(\cdot | x_k)$ for all $h_k \in \mathcal{H}_k$ and $k \in [N - 1]$. Otherwise, it is called causal. We denote the set of stochastic causal, stochastic Markov, deterministic causal and deterministic Markov policies by $\Pi, \Pi_{\text{m}}, \Pi_{\text{d}}$ and Π_{dm} , respectively, and note that $\Pi_{\text{dm}} \subseteq \Pi_{\text{m}} \subseteq \Pi$, and $\Pi_{\text{dm}} \subseteq \Pi_{\text{d}} \subseteq \Pi$. A mixed policy π_{mix} describes a deterministic policy that is randomly chosen at the initial time step $k = 0$ and then used during the entire control task. More formally, a mixed policy can be defined as in [25]. We endow the set of deterministic policies Π_{d} with a metric topology and define the corresponding Borel σ -algebra by $\mathcal{G}_{\Pi_{\text{d}}}$. A mixed policy π_{mix} can be seen as a random variable to the measurable space $(\Pi_{\text{d}}, \mathcal{G}_{\Pi_{\text{d}}})$. The set of all mixed policies is denoted as Π_{mix} .

For a given initial state $x_0 \in \mathcal{X}$, policy $\pi \in \Pi$, and the transition kernel T , a unique probability measure $\mathbb{P}_{x_0}^{\pi}$ for the state-input-trajectory is defined over $\mathcal{B}((\mathcal{X} \times \mathcal{U})^N \times \mathcal{X})$, which can be sampled recursively via $x_{k+1} \sim T(\cdot | x_k, u_k)$ with $u_k \sim \mu_k(h_k)$ [26]. We define the associated expected cumulative cost as

$$\mathbb{E}_{x_0}^{\pi} \left[\ell_N(x_N) + \sum_{k=0}^{N-1} \ell_k(x_k, u_k) \right] = \int_{(\mathcal{X} \times \mathcal{U})^N \times \mathcal{X}} \left(\ell_N(x_N) + \sum_{k=0}^{N-1} \ell_k(x_k, u_k) \right) \mathbb{P}_{x_0}^{\pi}(dx_0, du_0, \dots, dx_N)$$

and the probability of safety as

$$\mathbb{P}_{x_0}^{\pi}(x_{0:N} \in \mathcal{A}) = \int_{(\mathcal{X} \times \mathcal{U})^N \times \mathcal{X}} \prod_{k=0}^{N-1} \mathbb{1}_{\mathcal{A}}(x_k) \mathbb{P}_{x_0}^{\pi}(dx_0, du_0, \dots, dx_N),$$

respectively. Since the probability measure over trajectories $\mathbb{P}_{x_0}^{\pi}$ associated to any policy π and starting from the initial state x_0 is unique, the associated expected cost and safety is unique.

Our aim is to minimize the cumulative cost, while guaranteeing a prescribed level of safety over the entire duration of

the mission, expressed as

$$\begin{aligned} & \inf_{\pi \in \Pi} \mathbb{E}_{x_0}^{\pi} \left[\ell_N(x_N) + \sum_{k=0}^{N-1} \ell_k(x_k, u_k) \right] \\ & \text{subject to } \mathbb{P}_{x_0}^{\pi}(x_{0:N} \in \mathcal{A}) \geq \alpha, \end{aligned} \quad (1)$$

where $\alpha \in \mathbb{R}_{\geq 0}$ is a user-specified risk tolerance parameter. The constraint realizes the requirement for the state trajectory $x_{0:N}$ to lie in the safe set \mathcal{A} with a probability of at least α .

To ensure Problem (1) is well posed, we consider the following standing assumption throughout.

Assumption 1: The input set \mathcal{U} is compact. Furthermore, for every $x \in \mathcal{X}$ and $\mathcal{A} \in \mathcal{B}(\mathcal{X})$, the transition kernel $T(\mathcal{A}|x, u)$ and stage cost $\ell_k(x, u)$ are continuous with respect to u for all $k \in [N-1]$.

As shown later, if Problem (1) is feasible, the infimum is attained under Assumption 1. Note that, to streamline the presentation, we rely on assumptions that are more restrictive than needed; alternative sufficient conditions can be found in [27].

III. PRELIMINARIES

A. Dynamic Programming

In the absence of the chance constraint, for a given policy $\pi \in \Pi_{\text{dm}}$ and initial state $x_0 \in \mathcal{X}$, the total cost incurred

$$C_0^{\pi}(x_0) = \mathbb{E}_{x_0}^{\pi} \left[\ell_N(x_N) + \sum_{k=0}^{N-1} \ell_k(x_k, u_k) \right]$$

satisfies the DP recursion [28]

$$\begin{aligned} C_N^{\pi}(x_N) &= \ell_N(x_N), \\ C_k^{\pi}(x_k) &= \ell_k(x_k, u_k) + \int_{\mathcal{X}} C_{k+1}^{\pi}(x_{k+1}) T(dx_{k+1}|x_k, u_k) \end{aligned} \quad (2)$$

where $u_k \sim \mu_k(x_k)$. The finite-horizon optimal control problem now aims to find the infimum expected cost over the policies $\pi \in \Pi_{\text{dm}}$. Denoting $C_k^*(x_k) = \inf_{\pi \in \Pi_{\text{dm}}} C_k^{\pi}(x_k)$, it follows that [28]

$$\begin{aligned} C_N^*(x_N) &= \ell_N(x_N), \\ C_k^*(x_k) &= \inf_{u_k \in \mathcal{U}} \ell_k(x_k, u_k) + \int_{\mathcal{X}} C_{k+1}^*(x_{k+1}) T(dx_{k+1}|x_k, u_k), \end{aligned} \quad (3)$$

where the infimum is attained under Assumption 1. As the system dynamics are Markovian, the optimal deterministic Markov policy is also optimal within the class of stochastic causal policies [27, Theorem 3.2.1].

B. Safety via Dynamic Programming

Conversely, in the absence of a cost, the safety of a policy can also be encoded via a DP recursion with multiplicative cost. We use $V_k^{\pi}, V_k^* : \mathcal{X} \rightarrow [0, 1]$ as shorthand notation for $V_k^{\pi}(x_k) = \mathbb{P}(x_{k:N} \in \mathcal{A}|x_k, \pi)$ and $V_k^*(x_k) = \sup_{\pi \in \Pi_{\text{dm}}} \mathbb{P}(x_{k:N} \in \mathcal{A}|x_k, \pi)$, respectively. Following [26],

$$\begin{aligned} V_N^{\pi}(x_N) &= \mathbb{1}_{\mathcal{A}}(x_N), \\ V_k^{\pi}(x_k) &= \mathbb{1}_{\mathcal{A}}(x_k) \int_{\mathcal{X}} V_{k+1}^{\pi}(x_{k+1}) T(dx_{k+1}|x_k, u_k), \end{aligned} \quad (4)$$

where $u_k \sim \mu_k(x_k)$, and

$$\begin{aligned} V_N^*(x_N) &= \mathbb{1}_{\mathcal{A}}(x_N), \\ V_k^*(x_k) &= \sup_{u_k \in \mathcal{U}} \mathbb{1}_{\mathcal{A}}(x_k) \int_{\mathcal{X}} V_{k+1}^*(x_{k+1}) T(dx_{k+1}|x_k, u_k). \end{aligned} \quad (5)$$

Assumption 1 guarantees that the supremum in recursion (5) is attained [29].

C. Duality theory

Consider the optimization problem

$$\begin{aligned} f^* &= \inf_{x \in \mathbb{R}^n} f(x) \\ & \text{subject to } g_j(x) \leq 0, \quad j = 1, \dots, r, \end{aligned} \quad (6)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $g_j : \mathbb{R}^n \rightarrow \mathbb{R}$. We refer to (6) as the primal problem and we denote its value by f^* . The dual problem is given by

$$\begin{aligned} q^* &= \max_{\lambda_1, \dots, \lambda_r} q(\lambda) \\ & \text{subject to } \lambda_j \geq 0, \quad j = 1, \dots, r, \end{aligned} \quad (7)$$

where $q(\lambda) = \inf_{x \in \mathbb{R}^n} \{f(x) + \sum_{j=1}^r \lambda_j g_j(x)\}$ is called the dual function. The function $L(x, \lambda) = f(x) + \sum_{j=1}^r \lambda_j g_j(x)$ is referred to as the Lagrangian. The dual problem is always a convex optimization problem even if the primal is not convex [30]. In general, $q^* \geq f^*$; if $q^* = f^*$ we say that strong duality holds and there is no duality gap. For a given $\bar{\lambda} \in \mathbb{R}^r$, let $x_{\bar{\lambda}} \in \operatorname{argmin}_{x \in \mathbb{R}^n} L(x, \bar{\lambda})$ be a value minimizing the Lagrangian. Then, $g(x_{\bar{\lambda}})$ is a subgradient of the dual function q evaluated at $\bar{\lambda}$. The duality theory discussed above can be extended to infinite dimensional spaces, see [31].

IV. JOINT CHANCE CONSTRAINED DYNAMIC PROGRAMMING

We formulate Problem (1) to be solved over the class of stochastic causal policies as they are the most general. By interpreting the policy and transition kernel as strategies of two players in a game and invoking the results in [25], one can show that for any stochastic causal policy $\pi \in \Pi$, there exists a corresponding mixed policy $\pi_{\text{mix}} \in \Pi_{\text{mix}}$ that generates the same probability measure $\mathbb{P}_{x_0}^{\pi_{\text{mix}}} = \mathbb{P}_{x_0}^{\pi}$ for the state-input-trajectory, and vice-versa. Note that the cost and safety constraint in Problem (1) are fully defined by this probability measure. Therefore, w.l.o.g., we aim to solve Problem (1) over the class of mixed policies, which returns the same infimum value. In this section, we characterize the structure of the optimal mixed policy, which allows us to construct it using deterministic Markov policies, for which we solve using DP.

A. Lagrangian Dual Framework

In line with [4], we rely on the Lagrangian dual formulation of Problem (1), which, thanks to the equivalence of mixed and stochastic causal policies, is given by

$$\begin{aligned} & \max_{\lambda \in \mathbb{R}_{\geq 0}} \inf_{\pi \in \Pi_{\text{mix}}} \mathbb{E}_{x_0}^{\pi} \left[\ell_N(x_N) + \sum_{k=0}^{N-1} \ell_k(x_k, u_k) \right] \\ & \quad + \lambda (\alpha - \mathbb{P}_{x_0}^{\pi}(x_{0:N} \in \mathcal{A})). \end{aligned} \quad (8)$$

Note that the inner infimum can be equivalently cast as

$$\inf_{\pi \in \Pi_{\text{mix}}} \mathbb{E}_{x_0}^{\pi} \left[\ell_N(x_N) + \sum_{k=0}^{N-1} \ell_k(x_k, u_k) - \lambda \prod_{k=0}^N \mathbb{1}_{\mathcal{A}}(x_k) \right] + \lambda \alpha. \quad (9)$$

Unfortunately, this dual problem is notoriously hard to solve, even in the conceptually much simpler space of deterministic Markov policies. To see this, notice that one might naively try to apply DP by defining the terminal and stage costs as

$$\begin{aligned} \ell_{\lambda, N}(x_N) &= \ell_N(x_N) - \lambda \prod_{k=0}^N \mathbb{1}_{\mathcal{A}}(x_k) + \lambda \alpha, \\ \ell_{\lambda, k}(x_k, u_k) &= \ell_k(x_k, u_k). \end{aligned} \quad (10)$$

Note that the terminal cost depends on the full state trajectory, information we only have at time-step N , but whose distribution depends on the policy at time-steps $k = 0, \dots, N-1$. This non-Markovian structure prevents the use of DP.

To overcome this difficulty, we start with an intuitive observation: The product $\prod_{k=0}^N \mathbb{1}_{\mathcal{A}}(x_k)$ in the terminal cost function (10) depends on the state trajectory up to time N , but in the end it simply equals one if all states remained within the safe set \mathcal{A} , and zero otherwise. We can easily encode this information by introducing a binary state b_k that, inspired by the discussion in Section III-B, is initialized to $\mathbb{1}_{\mathcal{A}}(x_0)$, and is set to zero whenever the trajectory leaves the safe set. Consequently, we can formulate the product in (10) in terms of the value of the binary state at terminal time and recover the Markovian structure of the problem. More formally, we define as $\tilde{\mathcal{X}} = \mathcal{X} \times \{0, 1\}$ and $\tilde{x}_k = (x_k, b_k) \in \tilde{\mathcal{X}}$, with $b_0 = \mathbb{1}_{\mathcal{A}}(x_0)$ and $b_{k+1} = \mathbb{1}_{\mathcal{A}}(x_{k+1})b_k$, leading to the overall dynamics $\tilde{T} : \mathcal{B}(\tilde{\mathcal{X}}) \times \tilde{\mathcal{X}} \times \mathcal{U} \rightarrow [0, 1]$,

$$\begin{aligned} &\tilde{T}(\tilde{x}_{k+1} | \tilde{x}_k, u_k) \\ &= \begin{cases} T(x_{k+1} | x_k, u_k) & \text{if } (b_k = 0 \wedge b_{k+1} = 0) \\ & \vee (b_k = 1 \wedge b_{k+1} = \mathbb{1}_{\mathcal{A}}(x_{k+1})) \\ 0 & \text{otherwise,} \end{cases} \end{aligned}$$

or equivalently

$$\begin{aligned} \tilde{T}(\tilde{x}_{k+1} | \tilde{x}_k, u_k) &= T(x_{k+1} | x_k, u_k) \left((1 - b_k)(1 - b_{k+1}) \right. \\ &\quad \left. + b_k(1 - |b_{k+1} - \mathbb{1}_{\mathcal{A}}(x_{k+1})|) \right), \end{aligned}$$

which is again measurable. For readability, we overload the notation $\tilde{T}(\tilde{x}_{k+1} | \tilde{x}_k, u_k) = T(\tilde{x}_{k+1} | \tilde{x}_k, u_k)$, $\mathbb{1}_{\mathcal{A}}(\tilde{x}_k) = \mathbb{1}_{\mathcal{A}}(x_k)$ and $\ell_N(\tilde{x}_N) = \ell_N(x_N)$, $\ell_k(\tilde{x}_k, u_k) = \ell_k(x_k, u_k)$. Policies are defined as before replacing \mathcal{X} with $\tilde{\mathcal{X}}$.

The expected product of the set indicator functions in (9) is now given by the expected value of the binary state, $\mathbb{E}_{x_0}^{\pi} \left[\prod_{k=0}^N \mathbb{1}_{\mathcal{A}}(\tilde{x}_k) \right] = \mathbb{E}_{x_0}^{\pi} [b_N]$, and (9) becomes

$$\max_{\lambda \in \mathbb{R}_{\geq 0}} \inf_{\pi \in \Pi_{\text{mix}}} \mathbb{E}_{x_0}^{\pi} \left[\ell_N(\tilde{x}_N) + \sum_{k=0}^{N-1} \ell_k(\tilde{x}_k, u_k) + \lambda(\alpha - b_N) \right]. \quad (11)$$

The inner minimization problem is now a Markov Decision Problem whose optimal solution for a fixed λ is attained by

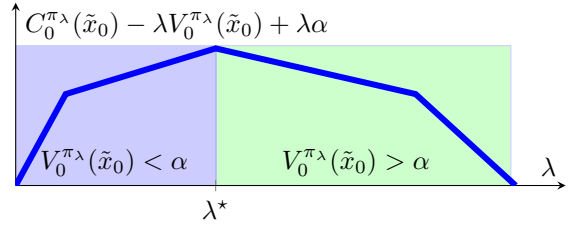


Fig. 2: An illustration of a objective function for the outer maximization in Problem (11) (inspired by [4]). The policy π_λ denotes an optimal argument to the inner minimization of Problem (11) under given λ (assuming it exists) and $C_0^{\pi_\lambda}(\tilde{x}_0)$, $V_0^{\pi_\lambda}(\tilde{x}_0)$ the cost and safety associated with that policy.

deterministic Markov policies and mixtures thereof, as we will show next.

B. Bilevel Framework

To simplify notation, note that using the equivalence $\mathbb{E}_{\tilde{x}_{1:N}} [b_N] = V_0^\pi(\tilde{x}_0)$, Problem (11) reduces to

$$\begin{aligned} &\max_{\lambda \in \mathbb{R}_{\geq 0}} \inf_{\pi \in \Pi_{\text{mix}}} \mathbb{E}_{x_0}^{\pi} \left[\ell_N(\tilde{x}_N) + \sum_{k=0}^{N-1} \ell_k(\tilde{x}_k, u_k) - \lambda b_N \right] + \lambda \alpha \\ &= \max_{\lambda \in \mathbb{R}_{\geq 0}} \inf_{\pi \in \Pi_{\text{mix}}} C_0^\pi(\tilde{x}_0) - \lambda V_0^\pi(\tilde{x}_0) + \lambda \alpha. \end{aligned} \quad (12)$$

Intuitively, this means that the inner Problem reduces to finding a policy that minimizes the expected cost while being rewarded by $-\lambda$ for remaining safe.

Referring to our introduction on duality in Section III-C, we note the following: We have the constraint $g(\pi) = \alpha - V_0^\pi(\tilde{x}_0)$ and Lagrangian

$$q(\lambda) = \inf_{\pi \in \Pi_{\text{mix}}} C_0^\pi(\tilde{x}_0) - \lambda V_0^\pi(\tilde{x}_0) + \lambda \alpha, \quad (13)$$

which is a concave function in λ and where we will show later that the infimum is attained by some $\pi_\lambda \in \Pi_{\text{mix}}$. Then $\alpha - V_0^{\pi_\lambda} \in \partial q(\lambda)$. Hence, if $V_0^{\pi_\lambda} < \alpha$ then $\lambda \leq \lambda^*$ and $\lambda \geq \lambda^*$ otherwise. A graphical illustration is given in Fig. 2.

Exploiting this structure, we propose to solve Problem (12) via the following auxiliary bilevel program

$$\begin{aligned} &\min_{\lambda \in \mathbb{R}_{\geq 0}} \lambda \\ &\text{subject to } \pi_\lambda \in \operatorname{argmin}_{\pi \in \Pi_{\text{mix}}} C_0^\pi(\tilde{x}_0) - \lambda V_0^\pi(\tilde{x}_0) \\ &\quad V_0^{\pi_\lambda}(\tilde{x}_0) \geq \alpha, \end{aligned} \quad (14)$$

where we drop the term $\lambda \alpha$ in the first constraint since it is a constant and does not affect the optimal argument π_λ .

To establish that (14) is equivalent to Problem (1), we show that, under mild assumptions, a mixed policy π_λ is indeed attainable for any $\lambda \in \mathbb{R}_{\geq 0}$, and that the policy π_{λ^*} associated with the optimal solution λ^* of (14) is the optimal solution to Problem (1). We divide the argument into three steps.

1) *Attainability of Optimal Deterministic Policies:* We first restrict ourselves to the space of deterministic policies and aim to find

$$\inf_{\pi \in \Pi_d} C_0^\pi(\tilde{x}_0) - \lambda V_0^\pi(\tilde{x}_0). \quad (15)$$

Since the problem is an MDP, it is sufficient to consider deterministic Markov policies, since no deterministic causal policy can perform better [27, Theorem 3.2.1]. In this case, the infimum can be easily computed using the recursion

$$\begin{aligned} J_N^*(\tilde{x}_N) &= \ell_N(\tilde{x}_N) - \lambda b_N, \\ J_k^*(\tilde{x}_k) &= \inf_{u_k \in \mathcal{U}} \ell_k(\tilde{x}_k, u_k) + \\ &\quad \int_{\tilde{\mathcal{X}}} J_{k+1}^*(\tilde{x}_{k+1}) T(d\tilde{x}_{k+1} | \tilde{x}_k, u_k). \end{aligned} \quad (16)$$

Theorem 4.1: (Attainability of Deterministic Markov Policies) Given a fixed $\lambda \in \mathbb{R}_{\geq 0}$, there exists a measurable deterministic Markov policy that attains the infimum in (16) at every time-step $k \in [N]$ and is also an optimal solution to Problem (15). Furthermore, the resulting $J_k^*(\cdot)$ is measurable for all $k \in [N]$.

Proof: Define $J_k^\pi(\tilde{x}_k) = C_k^\pi(\tilde{x}_k) - \lambda V_k^\pi(\tilde{x}_k) b_k$ and $J_k^*(\tilde{x}_k) = \inf_{\pi \in \Pi_{\text{dm}}} J_k^\pi(\tilde{x}_k)$, leading to $J_0^*(\tilde{x}_0) = \inf_{\pi \in \Pi_{\text{dm}}} C_0^\pi(\tilde{x}_0) - \lambda V_0^\pi(\tilde{x}_0) b_0 = \inf_{\pi \in \Pi_{\text{dm}}} C_0^\pi(\tilde{x}_0) - \lambda V_0^\pi(\tilde{x}_0)$ since $V_0^\pi(\tilde{x}_0)$ is zero if $x_0 \notin \mathcal{A}$, i.e., if b_0 is zero. Note that for any $\pi \in \Pi_{\text{dm}}$

$$\begin{aligned} V_k^\pi(\tilde{x}_k) b_k &= \mathbb{1}_{\mathcal{A}}(\tilde{x}_k) \int_{\tilde{\mathcal{X}}} V_{k+1}^\pi(\tilde{x}_{k+1}) b_k T(d\tilde{x}_{k+1} | \tilde{x}_k, \mu_k(\tilde{x}_k)) \\ &= \int_{\tilde{\mathcal{X}}} V_{k+1}^\pi(\tilde{x}_{k+1}) b_{k+1} T(d\tilde{x}_{k+1} | \tilde{x}_k, \mu_k(\tilde{x}_k)), \end{aligned}$$

and $V_N^\pi(\tilde{x}_N) b_N = \mathbb{1}_{\mathcal{A}}(\tilde{x}_N) \prod_{k=0}^N \mathbb{1}_{\mathcal{A}}(\tilde{x}_k) = b_N$. We have

$$\begin{aligned} J_N^*(\tilde{x}_N) &= \inf_{\pi \in \Pi_{\text{dm}}} C_N^\pi(\tilde{x}_N) - \lambda V_N^\pi(\tilde{x}_N) b_N \\ &= \ell_N(\tilde{x}_N) - \lambda b_N \\ J_k^*(\tilde{x}_k) &= \inf_{\pi \in \Pi_{\text{dm}}} C_k^\pi(\tilde{x}_k) - \lambda V_k^\pi(\tilde{x}_k) b_k \\ &= \inf_{\pi \in \Pi_{\text{dm}}} \ell_k(\tilde{x}_k, \mu_k(\tilde{x}_k)) + \int_{\tilde{\mathcal{X}}} \left(C_{k+1}^\pi(\tilde{x}_{k+1}) \right. \\ &\quad \left. - \lambda V_{k+1}^\pi(\tilde{x}_{k+1}) b_{k+1} \right) T(d\tilde{x}_{k+1} | \tilde{x}_k, \mu_k(\tilde{x}_k)) \\ &= \inf_{\pi \in \Pi_{\text{dm}}} \ell_k(\tilde{x}_k, \mu_k(\tilde{x}_k)) \\ &\quad + \int_{\tilde{\mathcal{X}}} J_{k+1}^\pi(\tilde{x}_{k+1}) T(d\tilde{x}_{k+1} | \tilde{x}_k, \mu_k(\tilde{x}_k)) \\ &= \inf_{u_k \in \mathcal{U}} \ell_k(\tilde{x}_k, u_k) + \int_{\tilde{\mathcal{X}}} J_{k+1}^*(\tilde{x}_{k+1}) T(d\tilde{x}_{k+1} | \tilde{x}_k, u_k) \end{aligned}$$

where the last equality follows since $T(\cdot | \tilde{x}_k, u_k) \geq 0$ and $J_{k+1}^*(\tilde{x}_{k+1})$ is the point-wise infimum of $J_{k+1}^\pi(\tilde{x}_{k+1})$ over $\pi \in \Pi_{\text{dm}}$ by definition and independent of u_k .

To complete the proof, we follow the induction argument in [29, Prop. 1]. Assume $J_{k+1}^*(\cdot)$ is measurable and let $F(\tilde{x}_k, u_k) = \ell_k(\tilde{x}_k, u_k) + \int_{\tilde{\mathcal{X}}} J_{k+1}^*(\tilde{x}_{k+1}) T(d\tilde{x}_{k+1} | \tilde{x}_k, u_k)$. Since $\ell_k(\tilde{x}_k, u_k)$ and $T(\tilde{x}_{k+1} | \tilde{x}_k, u_k)$ are continuous in u_k , $F(\tilde{x}_k, u_k)$ is continuous in u_k for every \tilde{x}_k [32, Fact 3.9]. Since \mathcal{U} is compact, the supremum is thus attained by a measurable map $\mu_k^*(\cdot)$ [33, Corollary 1]. Using the fact that

$F(\cdot, \cdot)$ and $\mu_k^*(\cdot)$ are measurable, we have that $F(\cdot, \mu_k^*(\cdot))$ is measurable; since $\ell_k(\cdot, \cdot)$ is measurable, $J_k^*(\cdot)$ is measurable [28, Prop. 7.29]. As $J_N^*(\cdot) = \ell_N(\cdot) - \lambda b_N$ is measurable, it follows by induction that $J_k^*(\cdot)$ is measurable and the infimum is attained by a measurable map $\mu_k^*(\cdot)$ for all $k \in [N-1]$. ■

2) *Attainability of Optimal Mixed Policies:* Throughout the remainder of the paper, we denote by $\bar{C}(\tilde{x}_0), \bar{V}(\tilde{x}_0)$ the control cost and safety associated with the maximum safety recursion (5) and $\underline{C}(\tilde{x}_0), \underline{V}(\tilde{x}_0)$ the control cost and safety associated with the minimum cost recursion (3).

Definition 4.2: For a given initial condition $\tilde{x}_0, \bar{C}(\tilde{x}_0)$ bounded, consider the set of performance attainable under any policy class Γ defined as

$$P_{\Gamma, \tilde{x}_0} = \bigcup_{\pi \in \Gamma} \{(C_0^\pi(\tilde{x}_0), V_0^\pi(\tilde{x}_0))\}$$

and the Pareto front $P_{\Gamma, \tilde{x}_0}^* : [\underline{V}(\tilde{x}_0), \bar{V}(\tilde{x}_0)] \rightarrow [\underline{C}(\tilde{x}_0), \bar{C}(\tilde{x}_0)]$, is

$$\begin{aligned} P_{\Gamma, \tilde{x}_0}^*(p) &= \inf_{C, V} C \\ &\text{subject to } (C, V) \in P_{\Gamma, \tilde{x}_0} \\ &\quad V \geq p. \end{aligned}$$

Further, given $\lambda \in \mathbb{R}$, we denote as λ -optimal under the policy class Γ any policy $\pi_\lambda \in \Gamma_\lambda$ with $\Gamma_\lambda = \{\pi' \in \Gamma : \pi' \in \text{argmin}_{\pi \in \Gamma} C_0^\pi(\tilde{x}_0) - \lambda V_0^\pi(\tilde{x}_0)\}$.

The Pareto front comprises the points in P_{Γ, \tilde{x}_0} , for which no cost reduction can be achieved without sacrificing safety and vice versa. Thus, it is monotone non-decreasing by definition (see also Fig. 3 for a graphical illustration).

Theorem 4.3: (Reduction to Mixed Policies) The Performance Set associated with mixed policies $P_{\Pi_{\text{mix}}, \tilde{x}_0}$ is the convex hull of the Performance Set P_{Π_d, \tilde{x}_0} of deterministic causal policies.

Proof: We first show that the performance of any mixed policy is contained in the convex hull of the Performance Set $\text{Conv}(P_{d, \tilde{x}_0})$ of deterministic causal policies. Recall that a mixed policy is defined as a mixture of deterministic causal policies. Let $\pi_1, \dots, \pi_M \in \Pi_d$, which are sampled with probabilities $\eta_1, \dots, \eta_M \in \mathbb{R}_{\geq 0}$, where $\sum_{i=0}^M \eta_i = 1$. Let $(C_1, V_1), \dots, (C_M, V_M) \in P_{d, \tilde{x}_0}$ be the pairs of control cost and safety associated with π_1, \dots, π_M . Then the corresponding mixed policy has safety $V = \sum_{i=0}^M \eta_i V_i$ and control cost $C = \sum_{i=0}^M \eta_i C_i$. Hence, $(C, V) \in \text{Conv}(P_{d, \tilde{x}_0})$.

We next show that any point in the convex hull $\text{Conv}(P_{d, \tilde{x}_0})$ of the Performance Set of Π_d is associated with a mixed policy. Let $(C, V) \in \text{Conv}(P_{d, \tilde{x}_0})$. Since $P_{d, \tilde{x}_0} \in \mathbb{R}^2$, there must exist two points $(C_1, V_1), (C_2, V_2) \in P_{d, \tilde{x}_0}$ associated to some $\pi_1, \pi_2 \in \Pi_d$, such that $C = \eta C_1 + (1 - \eta) C_2$ and $V = \eta V_1 + (1 - \eta) V_2, \eta \in [0, 1]$. Defining a policy $\pi_{\text{mix}} \in \Pi_{\text{mix}}$, consisting of $\pi_1, \pi_2 \in \Pi_d$, which are sampled with probabilities η and $1 - \eta$, respectively, we obtain a mixed policy with cost C and safety V . ■

Theorem 4.3 immediately leads to the following.

Corollary 4.4: The Pareto front $P_{\Pi_{\text{mix}}, \tilde{x}_0}^*(p)$ associated with mixed policies is convex.

Corollary 4.5: Any point in the Performance Set $P_{\Pi_{\text{mix}}, \tilde{x}_0}$ can be attained by mixing at maximum two deterministic causal policies.

The former follows by convexity of the Performance Set of mixed policies, whereas the latter has also been observed in [34] and [35]. In fact, [35] shows that the number of deterministic policies that need to be mixed for a constrained MDP is one plus the number of constraints, i.e., two policies for Problem (1).

Lemma 4.6 (Attainability of Mixed Policies): Fix $\lambda \in \mathbb{R}_{\geq 0}$. Then, there exists at least one measurable λ -optimal mixed policy $\pi_\lambda \in \operatorname{argmin}_{\pi \in \Pi_{\text{mix}}} C_0^\pi(\tilde{x}_0) - \lambda V_0^\pi(\tilde{x}_0)$ that equals a λ -optimal deterministic Markov policy. Furthermore, any mixture of λ -optimal deterministic Markov policies is a λ -optimal mixed policy.

Proof: By Corollary 4.5, any point in $P_{\Pi_{\text{mix}}, \tilde{x}_0}$ can be constructed using two deterministic causal policies $\pi_1, \pi_2 \in \Pi_{\text{d}}$ that are played with some probability $\eta \in [0, 1]$ and $\eta - 1$, respectively. Then,

$$\begin{aligned} & \inf_{\pi \in \Pi_{\text{mix}}} C_0^\pi(\tilde{x}_0) - \lambda V_0^\pi(\tilde{x}_0) \\ &= \inf_{\pi_1, \pi_2 \in \Pi_{\text{d}}} \eta (C_0^{\pi_1}(\tilde{x}_0) - \lambda V_0^{\pi_1}(\tilde{x}_0)) \\ & \quad + (1 - \eta) (C_0^{\pi_2}(\tilde{x}_0) - \lambda V_0^{\pi_2}(\tilde{x}_0)) \\ &= \inf_{\pi_1 \in \Pi_{\text{d}}} \eta (C_0^{\pi_1}(\tilde{x}_0) - \lambda V_0^{\pi_1}(\tilde{x}_0)) \\ & \quad + \inf_{\pi_2 \in \Pi_{\text{d}}} (1 - \eta) (C_0^{\pi_2}(\tilde{x}_0) - \lambda V_0^{\pi_2}(\tilde{x}_0)) \\ &= \inf_{\pi_1 \in \Pi_{\text{dm}}} \eta (C_0^{\pi_1}(\tilde{x}_0) - \lambda V_0^{\pi_1}(\tilde{x}_0)) \\ & \quad + \inf_{\pi_2 \in \Pi_{\text{dm}}} (1 - \eta) (C_0^{\pi_2}(\tilde{x}_0) - \lambda V_0^{\pi_2}(\tilde{x}_0)), \\ &= \min_{\pi_1 \in \Pi_{\text{dm}}} \eta (C_0^{\pi_1}(\tilde{x}_0) - \lambda V_0^{\pi_1}(\tilde{x}_0)) \\ & \quad + \min_{\pi_2 \in \Pi_{\text{dm}}} (1 - \eta) (C_0^{\pi_2}(\tilde{x}_0) - \lambda V_0^{\pi_2}(\tilde{x}_0)), \end{aligned}$$

where the last two equalities follow from the problem being an MDP, allowing us to restrict attention to deterministic Markov policies, and Theorem 4.1. If the minimizing deterministic Markov policy is unique, i.e., $\Pi_{\text{dm}, \lambda}$ is a singleton, then the mixed policy is unique and equals the deterministic Markov policy. If it is non-unique, also any mixture of λ -optimal deterministic Markov policies yields a λ -optimal mixed policy. Since the minimizing λ -optimal deterministic Markov policies are measurable (see Theorem 4.1), the mixed policy is measurable. ■

Corollary 4.7: Any λ -optimal mixed policy has zero probability of sampling a deterministic Markov policy that is not λ -optimal.

Proof: Let $\pi_1 \in \Pi_{\text{dm}, \lambda}$ and $\pi_2 \in \Pi_{\text{dm}} \setminus \Pi_{\text{dm}, \lambda}$, i.e., $C_0^{\pi_1}(\tilde{x}_0) - \lambda V_0^{\pi_1}(\tilde{x}_0) < C_0^{\pi_2}(\tilde{x}_0) - \lambda V_0^{\pi_2}(\tilde{x}_0)$. Then,

$$\begin{aligned} & \min_{\pi \in \Pi_{\text{mix}}} C_0^\pi(\tilde{x}_0) - \lambda V_0^\pi(\tilde{x}_0) \\ &= C_0^{\pi_1}(\tilde{x}_0) - \lambda V_0^{\pi_1}(\tilde{x}_0) \\ &< \eta (C_0^{\pi_1}(\tilde{x}_0) - \lambda V_0^{\pi_1}(\tilde{x}_0)) \\ & \quad + (1 - \eta) (C_0^{\pi_2}(\tilde{x}_0) - \lambda V_0^{\pi_2}(\tilde{x}_0)) \end{aligned}$$

whenever $\eta < 1$, i.e., if there is a non-zero probability of sampling π_2 . Hence, whenever there is a non-zero probability of sampling a deterministic Markov policy that is not λ -optimal, the resulting mixed policy can also not be λ -optimal since sampling any λ -optimal deterministic Markov policy with probability one yields a superior mixed policy. ■

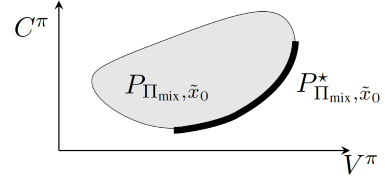


Fig. 3: Illustration of a Performance Set $P_{\Pi_{\text{mix}}, \tilde{x}_0}$ and its Pareto front $P_{\Pi_{\text{mix}}, \tilde{x}_0}^*$.

3) Equivalence of Problems:

Lemma 4.8: (Monotonicity) The safety $V_0^{\pi_\lambda}(\tilde{x}_0)$ and control cost $C_0^{\pi_\lambda}(\tilde{x}_0)$ of the policy $\pi_\lambda \in \Pi_{\text{mix}, \lambda}$ are monotone in $\lambda \in \mathbb{R}_{\geq 0}$. Specifically, if $0 \leq \lambda < \lambda'$, then $V_0^{\pi_\lambda}(\tilde{x}_0) \leq V_0^{\pi_{\lambda'}}(\tilde{x}_0)$ and $C_0^{\pi_\lambda}(\tilde{x}_0) \leq C_0^{\pi_{\lambda'}}(\tilde{x}_0)$ for all $\pi_\lambda \in \Pi_{\text{mix}, \lambda}, \pi_{\lambda'} \in \Pi_{\text{mix}, \lambda'}$.

Proof: Given some $\lambda > 0, \pi_\lambda \in \Pi_{\text{mix}, \lambda}$, we have that

$$C_0^{\pi_\lambda}(\tilde{x}_0) - \lambda V_0^{\pi_\lambda}(\tilde{x}_0) = \min_{\pi \in \Pi_{\text{mix}}} C_0^\pi(\tilde{x}_0) - \lambda V_0^\pi(\tilde{x}_0).$$

Hence, any policy with greater safety will be associated with greater control cost, and vice versa, i.e., for all $\pi \in \Pi_{\text{mix}}$

$$\begin{aligned} V_0^\pi(\tilde{x}_0) > V_0^{\pi_\lambda}(\tilde{x}_0) &\implies C_0^\pi(\tilde{x}_0) > C_0^{\pi_\lambda}(\tilde{x}_0), \\ C_0^\pi(\tilde{x}_0) < C_0^{\pi_\lambda}(\tilde{x}_0) &\implies V_0^\pi(\tilde{x}_0) < V_0^{\pi_\lambda}(\tilde{x}_0). \end{aligned}$$

Now consider $0 \leq \lambda < \lambda', \pi_{\lambda'} \in \Pi_{\text{mix}, \lambda'}$. Assume, for the sake of contradiction, that $V_0^{\pi_\lambda}(\tilde{x}_0) > V_0^{\pi_{\lambda'}}(\tilde{x}_0)$. Then, as above, $C_0^{\pi_\lambda}(\tilde{x}_0) > C_0^{\pi_{\lambda'}}(\tilde{x}_0)$ and note that $V_0^{\pi_\lambda}(\tilde{x}_0) > 0$ since $V_0^{\pi_\lambda}(\tilde{x}_0) > V_0^{\pi_{\lambda'}}(\tilde{x}_0) \geq 0$. We then have that

$$\begin{aligned} & C_0^{\pi_{\lambda'}}(\tilde{x}_0) - \lambda' V_0^{\pi_{\lambda'}}(\tilde{x}_0) \\ &\leq C_0^{\pi_\lambda}(\tilde{x}_0) - \lambda' V_0^{\pi_\lambda}(\tilde{x}_0) \\ &= C_0^{\pi_\lambda}(\tilde{x}_0) - \lambda V_0^{\pi_\lambda}(\tilde{x}_0) - (\lambda' - \lambda) V_0^{\pi_\lambda}(\tilde{x}_0) \\ &\leq C_0^{\pi_{\lambda'}}(\tilde{x}_0) - \lambda V_0^{\pi_{\lambda'}}(\tilde{x}_0) - (\lambda' - \lambda) V_0^{\pi_\lambda}(\tilde{x}_0) \\ &< C_0^{\pi_{\lambda'}}(\tilde{x}_0) - \lambda V_0^{\pi_{\lambda'}}(\tilde{x}_0) - (\lambda' - \lambda) V_0^{\pi_{\lambda'}}(\tilde{x}_0) \\ &= C_0^{\pi_{\lambda'}}(\tilde{x}_0) - \lambda' V_0^{\pi_{\lambda'}}(\tilde{x}_0), \end{aligned}$$

which is a contradiction.

The symmetric argument leads to the complementary inequalities, hence overall $V_0^{\pi_\lambda}(\tilde{x}_0) \leq V_0^{\pi_{\lambda'}}(\tilde{x}_0)$ and $C_0^{\pi_\lambda}(\tilde{x}_0) \leq C_0^{\pi_{\lambda'}}(\tilde{x}_0)$. Since this holds for any $\pi_\lambda \in \Pi_{\text{mix}, \lambda}, \pi_{\lambda'} \in \Pi_{\text{mix}, \lambda'}$, monotonicity follows. ■

Remark 4.9: A similar argument shows that Lemma 4.8 also holds for $\Pi_{\text{dm}, \lambda}$.

We are now ready to state the equivalence of Problems (1) and (14).

Theorem 4.10 (Equivalence of Problems (1) and (14)):

Assume that the minimum in (14) is attained by $\lambda^* \in [0, \infty)$. Then, there exists $\pi_{\lambda^*} \in \Pi_{\text{mix}, \lambda^*}$ such that $V_0^{\pi_{\lambda^*}}(\tilde{x}_0) \geq \alpha$ and π_{λ^*} is a minimizing argument of Problem (1).

Proof: We assume λ^* is attained, i.e., the constraints in Problem (14) must hold. Hence, there exists at least one λ^* -optimal policy $\pi_{\lambda^*} \in \Pi_{\text{mix}, \lambda^*}$ which yields $V_0^{\pi_{\lambda^*}}(\tilde{x}_0) \geq \alpha$.

Exploiting the equivalence of mixed and stochastic causal policies, let $\pi^* \in \Pi_{\text{mix}}$ be an optimal argument of Problem (1) with associated control cost $C_0^{\pi^*}(\tilde{x}_0)$ and safety $V_0^{\pi^*}(\tilde{x}_0)$.

Assume, for the sake of contradiction, that π_{λ^*} is not an optimal argument of Problem (1). Then, $C_0^{\pi_{\lambda^*}}(\tilde{x}_0) < C_0^{\pi^*}(\tilde{x}_0)$ by optimality of π^* in Problem (1) and consequently $V_0^{\pi_{\lambda^*}}(\tilde{x}_0) < V_0^{\pi^*}(\tilde{x}_0)$ by optimality of π_{λ^*} in Problem (14).

If there is a λ leading to $\pi^* \in \Pi_{\text{mix}, \lambda}$, we must have that for any $\bar{\pi}, \underline{\pi} \in \Pi_{\text{mix}}$ with $V_0^{\bar{\pi}}(\tilde{x}_0) \leq V_0^{\pi^*}(\tilde{x}_0) \leq V_0^{\underline{\pi}}(\tilde{x}_0)$,

$$\begin{aligned} C_0^{\pi^*}(\tilde{x}_0) - \lambda V_0^{\pi^*}(\tilde{x}_0) &\leq C_0^{\bar{\pi}}(\tilde{x}_0) - \lambda V_0^{\bar{\pi}}(\tilde{x}_0) \\ \Leftrightarrow \lambda V_0^{\bar{\pi}}(\tilde{x}_0) - \lambda V_0^{\pi^*}(\tilde{x}_0) &\leq C_0^{\bar{\pi}}(\tilde{x}_0) - C_0^{\pi^*}(\tilde{x}_0) \end{aligned}$$

and similarly

$$\begin{aligned} C_0^{\pi^*}(\tilde{x}_0) - \lambda V_0^{\pi^*}(\tilde{x}_0) &\leq C_0^{\underline{\pi}}(\tilde{x}_0) - \lambda V_0^{\underline{\pi}}(\tilde{x}_0) \\ \Leftrightarrow \lambda V_0^{\underline{\pi}}(\tilde{x}_0) - \lambda V_0^{\pi^*}(\tilde{x}_0) &\leq C_0^{\underline{\pi}}(\tilde{x}_0) - C_0^{\pi^*}(\tilde{x}_0), \end{aligned}$$

which combined leads to

$$\frac{C_0^{\bar{\pi}}(\tilde{x}_0) - C_0^{\pi^*}(\tilde{x}_0)}{V_0^{\bar{\pi}}(\tilde{x}_0) - V_0^{\pi^*}(\tilde{x}_0)} \leq \lambda \leq \frac{C_0^{\underline{\pi}}(\tilde{x}_0) - C_0^{\pi^*}(\tilde{x}_0)}{V_0^{\underline{\pi}}(\tilde{x}_0) - V_0^{\pi^*}(\tilde{x}_0)}.$$

Under these bounds a feasible λ always exists if the Pareto front is convex, which is indeed the case for mixed policies by Corollary 4.4.

However, then $\lambda < \lambda^*$ by Lemma 4.8. This violates the assumption that λ^* is the attained minimum. Hence π_{λ^*} is an optimal argument for Problem (1). ■

Remark 4.11: The equivalence of Problems (1) and (14) implies strong duality. This is highly surprising, given that it holds even when the cost functions and safe set are non-convex.

Interestingly, if the Pareto front $P_{\text{mix}, \tilde{x}_0}^*(\cdot)$ is strictly convex, then the optimal mixed policy π_{λ^*} will be equivalent to a Markov deterministic policy. The reason is that the set of λ^* -optimal Markov deterministic policies will become a singleton.

V. ALGORITHMIC SOLUTION

Next, we design an optimization procedure to solve Problem (14) and provide conditions to check feasibility.

A. Feasibility Check and Boundary Solutions

Following [4], one might first check the feasibility and triviality of Problem (1), where by the latter we mean that the minimum control cost policy is safe enough. Therefore, one might run recursions (3) and (5) and check whether the respective policies meet the desired safety level α . If the policy associated with the maximum safety recursion (5) yields a safety less than α , then Problem (1) is infeasible. If the policy associated with the minimum cost recursion (3) yields a safety of at least α , then it is optimal in terms of Problem (1) and there is no need to solve Problem (14).

Note that the policies associated with (3) and (5) are not necessarily unique and the recursions either only optimize for control cost or safety, respectively. However, it is still desirable to recover a preference for policies that, e.g., yield higher safety at the same control cost or lower control cost at the same safety. Therefore, the following Proposition is useful (see also Fig. 4).

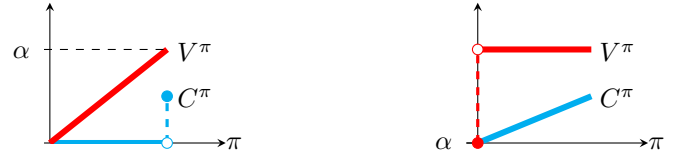


Fig. 4: In the left plot, we can choose from policies which have a safety arbitrarily close to α and a control cost of C or a policy that attains α but at cost $C + \delta$, $\delta > 0$. Then, for any λ , there always exists a policy π with safety close enough to α such that it is not optimal to incur the additional cost δ , i.e., $\lambda(\alpha - V^\pi) < \delta$. The right plot depicts a similar border case.

Proposition 5.1: (Border Case Suboptimality) Take any $\lambda > 0$. Then,

$$C^{\pi_\lambda}(\tilde{x}_0) - \underline{C}(\tilde{x}_0) \leq \lambda(V^{\pi_\lambda}(\tilde{x}_0) - \underline{V}(\tilde{x}_0)),$$

which goes to zero as λ goes to zero. Moreover, if $\bar{C}(\tilde{x}_0)$ is bounded. Then,

$$\bar{V}(\tilde{x}_0) - V^{\pi_\lambda}(\tilde{x}_0) \leq \frac{\bar{C}(\tilde{x}_0) - C^{\pi_\lambda}(\tilde{x}_0)}{\lambda},$$

which goes to zero as λ goes to infinity.

Proof: To show the first statement, fix $\lambda > 0$. By π_λ being λ -optimal, we have

$$C^{\pi_\lambda}(\tilde{x}_0) - \lambda V^{\pi_\lambda}(\tilde{x}_0) \leq \underline{C}(\tilde{x}_0) - \lambda \underline{V}(\tilde{x}_0).$$

Rearranging above equation yields the bound. Taking λ to zero, convergence to zero is guaranteed since $V^{\pi_\lambda}(\tilde{x}_0), \underline{V}(\tilde{x}_0)$ are bounded between zero and one. For the second statement, by π_λ being λ -optimal we have

$$C^{\pi_\lambda}(\tilde{x}_0) - \lambda V^{\pi_\lambda}(\tilde{x}_0) \leq \bar{C}(\tilde{x}_0) - \lambda \bar{V}(\tilde{x}_0)$$

Combined with $\bar{V}(\tilde{x}_0) \geq V^{\pi_\lambda}(\tilde{x}_0)$, we have $\bar{C}(\tilde{x}_0) \geq C^{\pi_\lambda}(\tilde{x}_0)$. Rearranging above equation yields the bound, which goes to zero with λ going to infinity since we assumed $\bar{C}(\tilde{x}_0)$ to be bounded and $\bar{C}(\tilde{x}_0) \geq C^{\pi_\lambda}(\tilde{x}_0) \geq 0$. ■

Most oftenly, however, the maximum safe policy is too costly and the minimum cost policy not safe enough, leading to a trade-off. We do not explore this trade-off much, but aim directly for a solution that attains the minimum cost at the required safety.

B. Bisection Algorithm

Let λ^* be the optimal solution to Problem (14) and recall that any associated λ^* -optimal mixed policy is constructed from a set of λ^* -optimal deterministic Markov policies. Consider $\pi_1, \pi_2 \in \Pi_{\text{dm}, \lambda^*}$, sampled with probabilities $\mu, 1 - \mu$, respectively. Then, the safety of the corresponding mixed policy π_{mix} is $V_0^{\pi_{\text{mix}}}(x_0) = \mu V_0^{\pi_1}(x_0) + (1 - \mu) V_0^{\pi_2}(x_0)$. For λ^* to be the optimal solution to Problem (14), we must have $V_0^{\pi_{\text{mix}}}(x_0) \geq \alpha$, and hence there must exist at least one policy $\pi \in \Pi_{\text{dm}, \lambda^*}$, which has a safety of at least α .

Now, by optimality of Problem (14), for any $\underline{\lambda} \in \mathbb{R}_{\geq 0}$ with $\underline{\lambda} < \lambda^*$, we have that $V_0^{\pi_{\text{mix}}}(x_0) < \alpha$. This implies that there does not exist $\pi \in \Pi_{\text{dm}, \underline{\lambda}}$ with $V_0^\pi(x_0) \geq \alpha$. Otherwise, choosing π_{mix} to sample π with probability one would yield

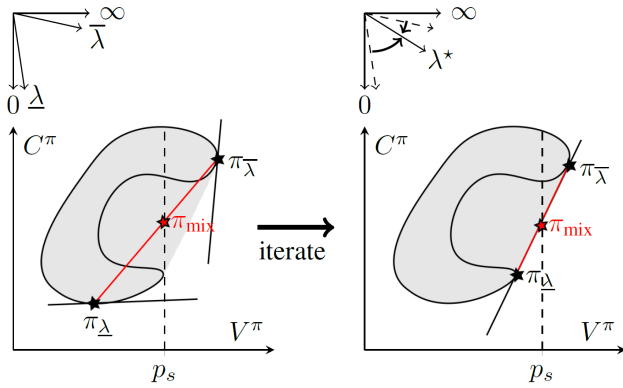


Fig. 5: Performance sets P_{Π_d, \tilde{x}_0} (bordered set) and its convex hull $P_{\Pi_{mix}, \tilde{x}_0}$ (grey set), as well as the performance of the respective policies $\underline{\lambda}, \bar{\lambda}$ (black stars) and the optimal interpolation π_{mix} according to equation (17) (red star). The variable λ defines the optimization direction. The DP recursion returns the optimal policy in the performance set in this direction. Along the outer loop iterations $\underline{\lambda}$ approaches $\bar{\lambda}$.

a feasible mixed policy. On the other hand, for any $\bar{\lambda}$ with $\bar{\lambda} > \lambda^*$, every $\bar{\lambda}$ -optimal deterministic Markov policy has a safety of at least α , since there exists a policy in Π_{dm, λ^*} with a safety of at least α (Remark 4.9).

Hence, based on two policies $\pi_{\underline{\lambda}} \in \Pi_{dm, \underline{\lambda}}$ and $\pi_{\bar{\lambda}} \in \Pi_{dm, \bar{\lambda}}$, we can construct a mixed policy with safety greater or equal to α , as assigning probability one to $\pi_{\bar{\lambda}}$ and zero to $\pi_{\underline{\lambda}}$ is always feasible. Ideally, one would assign an as high as possible probability measure to $\pi_{\underline{\lambda}}$, since it is associated with a lower control cost (Remark 4.9). This is achieved by linearly interpolating between $\pi_{\underline{\lambda}}$ and $\pi_{\bar{\lambda}}$ in such a way, that, in expectation, we obtain a safety of exactly $\alpha = p_{\bar{\lambda}} V_0^{\pi_{\bar{\lambda}}}(\tilde{x}_0) + (1 - p_{\bar{\lambda}}) V_0^{\pi_{\underline{\lambda}}}(\tilde{x}_0)$ (see Fig. 5). This in turn can be achieved by the policy

$$\pi_{mix} = \begin{cases} \pi_{\bar{\lambda}} & \text{with probability } p_{\bar{\lambda}} = \frac{\alpha - V_0^{\pi_{\underline{\lambda}}}(\tilde{x}_0)}{V_0^{\pi_{\bar{\lambda}}}(\tilde{x}_0) - V_0^{\pi_{\underline{\lambda}}}(\tilde{x}_0)}, \\ \pi_{\underline{\lambda}} & \text{otherwise.} \end{cases} \quad (17)$$

Following [4], we propose a bisection algorithm to shrink bounds $[\underline{\lambda}, \bar{\lambda}]$ on λ^* . Recursions (3) and (5) allow us to provide an initial range $[\underline{\lambda}_{init}, \bar{\lambda}_{init}]$ of values for λ^* .

Proposition 5.2 (Bounds on λ^* , [18]): The value $\underline{\lambda}_{init} = 0$ always yields a lower bound on $\lambda^* \in \mathbb{R}_{\geq 0}$. Assume $\bar{C}(\tilde{x}_0)$ is bounded and $\bar{V}(\tilde{x}_0) > \alpha$ and choose $\bar{\lambda}_{init} = \frac{\bar{C}(\tilde{x}_0) - \underline{C}(\tilde{x}_0)}{\bar{V}(\tilde{x}_0) - \alpha}$. Then, the policy $\pi_{\bar{\lambda}_{init}}$ achieves a safety greater than α .

Note that in [18] this guarantee is given for deterministic Markov policies. However, since any $\bar{\lambda}_{init}$ -optimal mixed policy can be constructed by $\bar{\lambda}_{init}$ -optimal deterministic Markov policies, $\bar{\lambda}_{init}$ is a feasible, yet not necessarily optimal solution to Problem (14).

Starting with the range provided in Proposition 5.2 we can now solve for λ^* using bisection, shrinking the bounds $\underline{\lambda} \leq \lambda^* \leq \bar{\lambda}$. The bisection algorithm picks λ in the middle of the interval $[\underline{\lambda}, \bar{\lambda}]$, then sets λ as the upper bound if it is greater than λ^* , and the lower bound otherwise, halving the interval. The test whether λ is greater or lower than λ^* is

performed by evaluating $V_0^{\pi_{\lambda}}$ and comparing it with α (see Lemma 4.8 and recall Fig. 2). Mixing the associated policies $\pi_{\underline{\lambda}}$ and $\pi_{\bar{\lambda}}$ as described in (17) yields a mixed policy, whose suboptimality to the solution of Problem (1) converges to zero at least exponentially over the number of bisection steps.

Theorem 5.3 (Suboptimality Bound, [18]): Let the bisection algorithm be initialized with $[\underline{\lambda}_{init}, \bar{\lambda}_{init}]$. After M bisection steps, the suboptimality of the policy constructed in (17) compared to the solution $\pi^* \in \Pi$ of Problem (1) is bounded by

$$\begin{aligned} C_0^{\pi_{mix}}(\tilde{x}_0) - C_0^{\pi^*}(\tilde{x}_0) &\leq p_{\bar{\lambda}}(1 - p_{\bar{\lambda}})(\bar{\lambda} - \underline{\lambda})(V_0^{\pi_{\bar{\lambda}}}(\tilde{x}_0) - V_0^{\pi_{\underline{\lambda}}}(\tilde{x}_0)) \\ &\leq 0.25 \left(\frac{1}{2}\right)^M (\bar{\lambda}_{init} - \underline{\lambda}_{init}). \end{aligned}$$

The process is summarized in Algorithm 1, which runs the bisection algorithm until a prescribed suboptimality gap Δ is attained, which is guaranteed to be achieved in finitely many iterations by Theorem 5.3. The computational performance is thus mainly determined by the inner loop, which is solved using DP. For continuous state and action spaces, the value function has to be approximated, e.g., by using gridding or basis functions [28], [36], [37]. Note that our results trivially extend to the case of finite state or input spaces.

Algorithm 1: Joint Chance Constr. Optimal Control

Data: $\mathcal{X}, \mathcal{U}, T, N, \tilde{x}_0, \alpha, \Delta$

```

/* Check border cases */
 $\bar{V}_0(\tilde{x}_0) \leftarrow$  Maximum safety recursion (5);
 $\underline{V}_0(\tilde{x}_0) \leftarrow$  Safety of minimum cost policy, (3) and (4);
if  $\bar{V}_0(\tilde{x}_0) \leq \alpha$  or  $\underline{V}_0(\tilde{x}_0) \geq \alpha$  then
  | Use method in Proposition 5.1
end

/* Run recursion */
 $\underline{\lambda} \leftarrow 0, \lambda \leftarrow \bar{\lambda} \leftarrow$  Proposition 5.2 ;
while true do
   $p_{\bar{\lambda}} \leftarrow$  Equation (17) ;
   $\delta \leftarrow p_{\bar{\lambda}}(1 - p_{\bar{\lambda}})(\bar{\lambda} - \underline{\lambda})(V_0^{\pi_{\bar{\lambda}}}(\tilde{x}_0) - V_0^{\pi_{\underline{\lambda}}}(\tilde{x}_0))$  ;
  if  $\delta \leq \Delta$  then
    | return  $p_{\bar{\lambda}}, \pi_{\underline{\lambda}}, \pi_{\bar{\lambda}}$  ;
  end
   $V_0^{\pi_{\lambda}}(\tilde{x}_0), \pi_{\lambda} \leftarrow$  Theorem 4.1 ;
  if  $V_0^{\pi_{\lambda}}(\tilde{x}_0) \leq \alpha$  then
    |  $\underline{\lambda} \leftarrow \lambda$  ;
  else
    |  $\bar{\lambda} \leftarrow \lambda$  ;
  end
   $\lambda \leftarrow \frac{1}{2}(\bar{\lambda} + \underline{\lambda})$  ;
end

```

VI. NUMERICAL EXAMPLE

We compare our algorithmic solution to the one proposed by [4] that relies on Boole's inequality to break the time correlation introduced by the joint chance constraints. Following

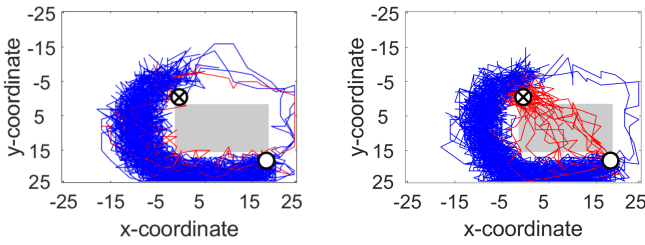


Fig. 6: Trajectories of 150 Monte Carlo simulations of the quadcopter in Example 1. The white area denotes the safe set \mathcal{A} , while the grey region should be avoided. The black circle denotes the initial state, the black circle with the cross the origin, which serves as the target state with zero stage cost. Blue curves denote safe trajectories, while red curves denote unsafe trajectories. The left plot shows the results for the policy of [4], the right plot the policy obtained by the recursion in Theorem 4.1. Despite the different behaviour, both policies provide the same safety probability of 0.9; the one on the right, however, attains a lower average cost.

[4], the application of Boole's inequality $\mathbb{P}_{x_0}^{\pi}(x_{0:N} \in \mathcal{A}) \geq 1 - \sum_{k=0}^{N-1} \mathbb{P}_{x_0}^{\pi}(x_k \in \mathcal{A}^c)$ to Problem (1) leads to the following Lagrangian dual

$$\begin{aligned} \max_{\lambda \in \mathbb{R}_{\geq 0}} \inf_{\pi \in \Pi_{\text{dm}}} \mathbb{E}_{x_0}^{\pi} \left[\ell_N(x_N) + \sum_{k=0}^{N-1} \ell_k(x_k, u_k) \right] \\ + \lambda \left(\sum_{k=0}^{N-1} \mathbb{P}(x_k \in \mathcal{A}^c) - 1 + \alpha \right), \end{aligned} \quad (18)$$

which displays a stage-wise decomposition structure. This allows the direct application of Dynamic Programming to compute the infimum, by using the terminal and stage-cost

$$\begin{aligned} \ell_{\lambda,N}(x_N) &= \ell_N(x_N) + \lambda \mathbb{1}_{\mathcal{A}^c}(x_N), \\ \ell_{\lambda,k}(x_k, u_k) &= \ell_k(x_k, u_k) + \lambda \mathbb{1}_{\mathcal{A}^c}(x_k). \end{aligned} \quad (19)$$

We compare this approach to the one proposed here via two examples on quadcopter trajectory planning and one example on fishery management. The code used to generate the results is available at https://github.com/NiklasSchmidResearch/JCC_opt_control.git.

Example 1 (Quadcopter Trajectory Planning I): We consider a quadcopter with state $x_k \in \mathcal{X} = [-25, 25] \times [-25, 25]$ that would like to minimize its cost defined as $\ell_k(x_k, u_k) = x_k^\top x_k$, for $k \in [20]$, while staying in the set \mathcal{A} with a predefined probability. We consider an Euler discretized unicycle model with constant speed of 3, given by $x_{k+1} = x_k + 3 [\cos(u) \quad \sin(u)]^\top + w_k$, where $w_k = \mathcal{N}([0 \ 0]^\top, \text{diag}(5, 5))$ is an additive disturbance, $u \in [0, 2\pi]$ is the input, and $x_0 = [13 \quad 13]^\top$.

We grid the space space into 50×50 discrete states and 8 discrete inputs to compute two policies: one by solving (19) with $\lambda = 8205$ and the other by solving (16) with $\lambda = 12645$, where the respective λ -values have been obtained by running Algorithm 1 with the respective DP recursions (19) and (16); the resulting policies are shown in Fig. 7. Over 150 Monte Carlo runs (Fig. 6) both policies result in a safety probability

of 0.9, however the policy developed in this work achieves a lower average cost of 7260 vs. 7580. This is because the policy from [4] tends to leave the restricted area as fast as possible after entering it. On the other hand, our policy keeps transitioning through the restricted area once it enters, since as soon as the trajectory is tagged as unsafe the policy only aims to minimize the expected cost. This difference is due to the structure of the DP recursions employed by the two algorithms. In our DP recursion (see Theorem 4.1) we enforce a one-time penalty of λ the first time the trajectory runs unsafe but not thereafter, while in (19) the policy incurs a penalty of λ for every time step that the state is outside the safe set.

Example 2: (Quadcopter Trajectory Planning II) The setting is adapted from Example 1, with all parameters remaining the same, except for the unsafe set that is now placed in the middle of the state space so that the cheapest states are unsafe and the initial state being $x_0 = [19 \quad 19]^\top$ (see the left plot in Fig. 8). To compare our DP recursion in Theorem 4.1 with that in [4] we evaluate the Pareto fronts of the two recursions by varying $\lambda \in [100, 10^6]$ (see the right plot in Fig. 8). We compare

- **(Dashed):** Evaluating the policy via Theorem 4.1 and its associated safety via (4) (our approach).
- **(Dash-Dotted):** Evaluating the policy via (19) and its associated safety via (4). Compared to the dashed graph, this shows the conservatism introduced by Boole's inequality in the DP recursion (19).
- **(Dotted)** Evaluating the policy via (19) and its associated safety via Boole's inequality [4]. Comparing the dotted with the dash-dotted graph yields the conservatism of Boole's inequality for evaluating the safety of a policy. Comparing the dotted with the dashed graph yields the performance difference between the approach by [4] and ours.

Clearly, the dashed graph dominates the dash-dotted graph, which in turn dominates the dotted graph, showing the effectiveness of our approach in reducing conservatism. For $\lambda = 0$, all methods are expected to generate policies that achieve the minimum achievable cost, as constraint violations are not penalized. This can be seen by the dashed and dash-dotted graph converging for small probabilities (the dotted plot achieves this cost at negative safety probabilities).

Interestingly, the Pareto fronts might not necessarily converge for high values of λ . Consider the fictitious example with two policies π_1, π_2 , where π_1 is the safest policy and has a probability of 0.1 to become unsafe at the first time-step and zero otherwise, and π_2 has a probability of 0.2 to be unsafe at the very last time-step and zero for all other time-steps. Let $N = 3$. Even if $C_0^{\pi_1}(\tilde{x}_0) = C_0^{\pi_2}(\tilde{x}_0)$, the DP recursion in [4] will prefer policy π_2 for any $\lambda > 0$ since the incurred penalty will be 0.3λ and 0.2λ for policy π_1 and π_2 , respectively. Thus, although π_1 might be the safest policy, it will never be considered optimal in [4], even for arbitrarily large λ . In our example, this can be observed by the graphs approaching, but never really touching each other for higher safety values. The only exception is when a safety of one is attainable since the conservatism of Boole's inequality converges to zero as

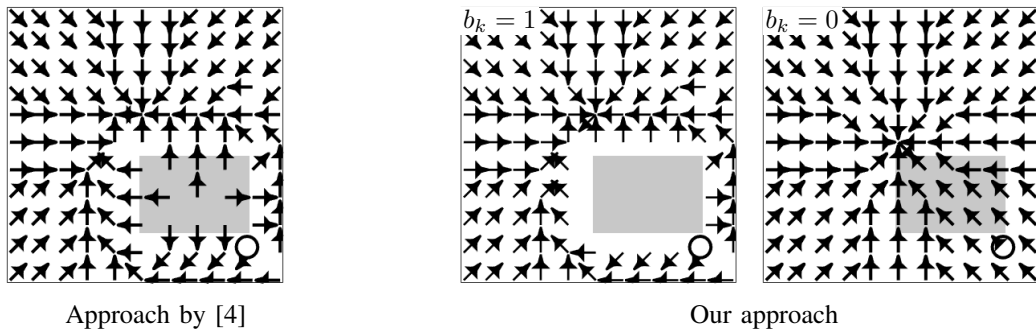


Fig. 7: Optimal inputs at time-step $k = 0$ associated to the recursion in (19) (left) and $\pi_{\bar{\lambda}}$ (middle for $b_k = 1$ and right for $b_k = 0$) for Example 1. Safe states are marked white, unsafe states grey. The initial state is marked with a black circle.

the safety approaches one. However, if we aim for lower requirements on safety, or if the highest achievable safety is far enough from one, the conservatism becomes significant. This can especially be observed for small safety values, e.g., to guarantee a safety of 0.1, the approach by [4] proposes a policy with cost 4113, while the DP recursion (19) is able to compute a policy with cost 3330 and our algorithm finds a policy with cost 1522.

We also tested Algorithm 1 using a desired safety of $\alpha = 0.6$ and a suboptimality bound of $\Delta = 10^{-6}$. The desired optimality was achieved after 20 iterations. For the approach in [4] we ran a fixed number of 30 iterations. Our approach returned a policy with cost 4122 compared to 5822 using the approach in [4].

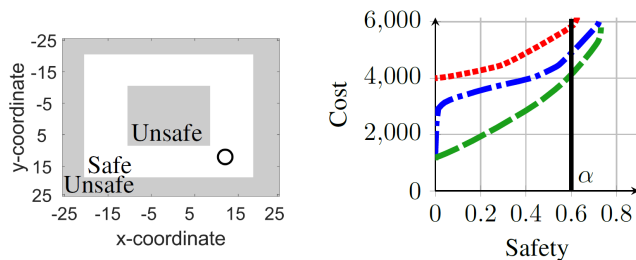


Fig. 8: The left plot shows the safe (white) and unsafe (grey) regions of the state space and the initial state (black circle) in example 2. The right plot shows the Pareto fronts when using our algorithm (Dashed), our algorithm but replacing the DP recursion with (19) (Dash-dotted), and when using the full approach in [4] (Dotted). In fact, the dotted graph continuous to negative probabilities due to conservativeness of Boole's inequality.

Example 3 (Fisheries Management): We adapt the fisheries management example from [38], which uses the model from [39]. The evolution of fish biomass in a reservoir x_k is described by

$$x_{k+1} = (1 - v_k)x_k + \gamma_k R(x_k) - C(x_k, u_k),$$

where $R(\cdot)$, $C(\cdot, \cdot)$ are functions representing recruitment and catch, $v_k \sim \mathcal{N}(0.2, 0.01)$ the natural mortality rate and $\gamma_k \sim \mathcal{N}(1, 0.36)$ variability in the recruitment. The catch function

is imposed by the government and described by

$$C(x_k, u_k) = \max \left\{ \delta_k u_k C, \delta_k u_k C \frac{x_k}{L} \right\},$$

where $\delta_k \sim \mathcal{N}(1.1, 0.04)$ is a variability in the catch with $\delta_k > 0$, $u_k \in [0, 1]$ our input variable denoting catch effort, C the maximum catch, and L the biomass limit of the reservoir. The recruitment function is described by

$$R(x_k) = x_k \left(1 - \frac{x_k}{L}\right) \text{sgm} \left(\frac{x_k - \mu}{\sigma^2} \right),$$

where the term $(1 - x_k/L)$ models decline in recruitment with saturation of the reservoir. The sigmoid $\text{sgm}(x) = \frac{1}{1 + e^{-x}}$ has been added here to introduce a bifurcation in the dynamics and could model a rapid decline in recruitment when the population becomes too small.

We use $L = 40$, $M = 10$, $\mu = 20$, $\sigma = 5$. Empirically, using Monte Carlo simulations, we found that the fish population is at almost zero probability of recovering whenever the biomass gets below 13 units, even if we stop fishing. Our goal is to catch as much fish as possible within 100 time-steps, while guaranteeing a probability of at least $\alpha = 0.75$ to preserve at least 13 units of biomass in the reservoir throughout. We define $\mathcal{X} = [0, 60]$ and discretize the state space into 60 discrete states, the input space $\mathcal{U} = [0, 1]$ into 5 discrete inputs and use $\Delta = 10^{-6}$.

Due to the conservatism of Boole's inequality, a solution to this problem would be infeasible in [4]. However, to allow for a comparison with our results, we again run our algorithm, replacing our DP recursion with equation (19). The policies computed this way and with our method are depicted in Fig. 9.

Intuitively, the less fish remains, the less aggressive the catch effort. Interestingly, however, our approach tends to fish more aggressively early on and the policy generated by [4] tends to fish more aggressively towards the end. The reason is that both methods aim for maximum catch with a certain risk of a population crash. However, while we penalize violation of constraints once, the penalty is repeated and accumulated in recursion (19) used by [4]. Hence, our method tries to exploit situations in which the population has high risk of crashing and tries to fish whatever remains, while the approach by [4] waits until the end to avoid accumulating a penalty of λ many times and then purposely crashes the population to maximize

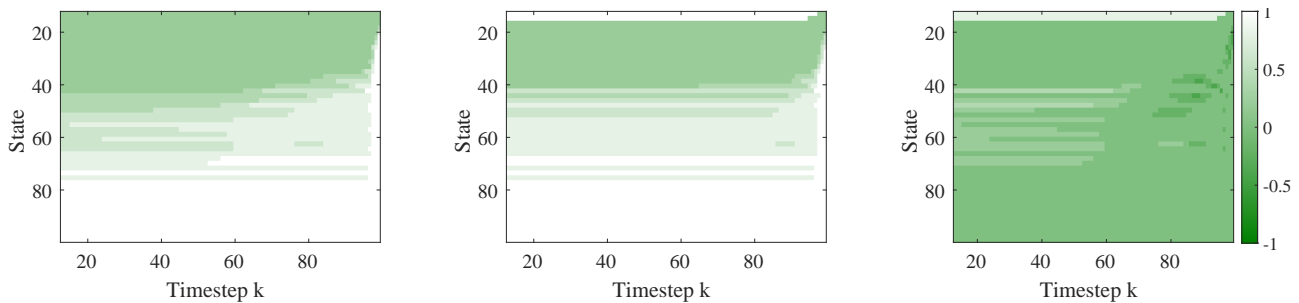


Fig. 9: Optimal policies for our fisheries management Example 3 using the approach in [4] (left), our approach (middle, showing $\pi_{\bar{\lambda}}$ for the case $b_k = 1$), and their difference (right plot). In the left two plots, the lighter the color, the higher the catch effort. In the right plot, the higher the catch effort difference the lighter the color, where positive values mean that our approach advertises a higher catch effort than that in [4].

catch. Overall, our method allows for a higher catch of 112.9 units compared to 102.3 units at a safety of 75%.

In comparison with [4] our algorithm tends to be computationally slightly more efficient. The reason is that the policy associated with the minimum control cost in (3) is applied until the end of the time horizon whenever b_k is zero, which is independent of λ and precomputed by the border case check. Thus, when λ is updated via bisection, it suffices run the DP recursions over states which have $b_k = 1$, which only exist within the safe set \mathcal{A} . In [4], on the other hand, the policy has to be reevaluated for all states in the state space \mathcal{X} as λ changes.

VII. CONCLUSIONS AND FUTURE WORKS

We consider the problem of optimally controlling stochastic systems subject to joint chance constraints over a finite-time horizon and proposed a Dynamic Programming scheme to solve for the respective optimal policy. Our analysis reveals interesting insights about this class of problems and uncovers a behaviour of the optimal policy induced by Problem (1), that may be controversial for many applications. Indeed, in practice one would avoid to play actions that are known to likely lead to constraint violations, no matter the cost. However, Problem (1) just considers safety and cost in expectation, meaning over infinitely many trials starting from the initial state. This is in line with our introduction of a binary state and mixed policies, which deliberately fail on constraints to achieve a cost reduction in expectation. On the other hand, trying to avoid deliberative failure by adding a constant penalty for unsafe states, as implicitly done in [4], is not solving Problem (1) to optimality and, furthermore, not necessarily generating a meaningful behaviour either. The desired behaviour of a policy within the unsafe region of the state space is highly application dependent and often requires another problem formulation than Problem (1). Specifically, Problem (1) is tailored to applications which reset the system state to the initial state after N time-steps, played infinitely often, and where constraints become redundant once they are violated, e.g., because an unrecoverable event is triggered. The formulation of an alternative problem that fits larger classes of application scenarios remains an open question and will be the focus of future work.

- [1] M. Prandini and J. Hu, "Application of reachability analysis for stochastic hybrid systems to aircraft conflict prediction," in *2008 47th IEEE conference on decision and control*. IEEE, 2008, pp. 4036–4041.
- [2] K. Lesser, M. Oishi, and R. S. Erwin, "Stochastic reachability for control of spacecraft relative motion," in *52nd IEEE Conference on Decision and Control*. IEEE, 2013, pp. 4705–4712.
- [3] M. Bojarski, D. Del Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang *et al.*, "End to end learning for self-driving cars," *arXiv preprint arXiv:1604.07316*, 2016.
- [4] Y. K. M. Ono, M. Pavone and J. Balam, "Chance-constrained dynamic programming with application to risk-aware robotic space exploration," *Autonomous Robots*, 2015.
- [5] G. Männel, J. Graßhoff, P. Rostalski, and H. S. Abbas, "Iterative gaussian process model predictive control with application to physiological control systems," in *2021 60th IEEE Conference on Decision and Control (CDC)*. IEEE, 2021, pp. 2203–2210.
- [6] M. Farina, L. Giulioni, and R. Scattolini, "Stochastic linear model predictive control with chance constraints—a review," *Journal of Process Control*, vol. 44, pp. 53–67, 2016.
- [7] M. Lorenzen, F. Dabbene, R. Tempo, and F. Allgöwer, "Constraint-tightening and stability in stochastic model predictive control," *IEEE Transactions on Automatic Control*, vol. 62, no. 7, pp. 3165–3177, 2016.
- [8] J. Köhler, F. Geuss, and M. N. Zeilinger, "On stochastic mpc formulations with closed-loop guarantees: Analysis and a unifying framework," *arXiv preprint arXiv:2304.00069*, 2023.
- [9] N. Schmid, J. Gruner, H. S. Abbas, and P. Rostalski, "A real-time gp based mpc for quadcopters with unknown disturbances," in *2022 American Control Conference (ACC)*. IEEE, 2022, pp. 2051–2056.
- [10] K. Wang and S. Gros, "Solving mission-wide chance-constrained optimal control using dynamic programming," in *2022 IEEE 61st Conference on Decision and Control (CDC)*. IEEE, 2022, pp. 2947–2952.
- [11] V. Raghuraman and J. P. Koeln, "Long duration stochastic mpc with mission-wide probabilistic constraints using waysets," *IEEE Control Systems Letters*, vol. 7, pp. 865–870, 2022.
- [12] M. Ono, "Joint chance-constrained model predictive control with probabilistic resolvability," in *2012 American Control Conference (ACC)*. IEEE, 2012, pp. 435–441.
- [13] J. A. Paulson, E. A. Buehler, R. D. Braatz, and A. Mesbah, "Stochastic model predictive control with joint chance constraints," *International Journal of Control*, vol. 93, no. 1, pp. 126–139, 2020.
- [14] K. Wang and S. Gros, "Recursive feasibility of stochastic model predictive control with mission-wide probabilistic constraints," in *2021 60th IEEE Conference on Decision and Control (CDC)*. IEEE, 2021, pp. 2312–2317.
- [15] A. Thorpe, T. Lew, M. Oishi, and M. Pavone, "Data-driven chance constrained control using kernel distribution embeddings," in *Learning for Dynamics and Control Conference*. PMLR, 2022, pp. 790–802.
- [16] V. A. Bavdekar and A. Mesbah, "Stochastic nonlinear model predictive control with joint chance constraints," *IFAC-PapersOnLine*, vol. 49, no. 18, pp. 270–275, 2016.
- [17] A. Charnes and W. W. Cooper, "Chance-constrained programming," *Management Science*, vol. 6, no. 1, pp. 73–79, 1959.

- [18] L. Pfeiffer, "Two approaches to constrained stochastic optimal control problems," *SFB Report*, vol. 7, p. 2015, 2015.
- [19] A. Patil and T. Tanaka, "Upper and lower bounds for end-to-end risks in stochastic robot navigation," *IFAC-PapersOnLine*, vol. 56, no. 2, pp. 5603–5608, 2023.
- [20] W. Chen, D. Subramanian, and S. Paternain, "Policy gradients for probabilistic constrained reinforcement learning," in *2023 57th Annual Conference on Information Sciences and Systems (CISS)*. IEEE, 2023, pp. 1–6.
- [21] M. Xu, Z. Liu, P. Huang, W. Ding, Z. Cen, B. Li, and D. Zhao, "Trustworthy reinforcement learning against intrinsic vulnerabilities: Robustness, safety, and generalizability," *arXiv preprint arXiv:2209.08025*, 2022.
- [22] Y. Zhang, Q. Vuong, and K. Ross, "First order constrained optimization in policy space," *Advances in Neural Information Processing Systems*, vol. 33, pp. 15 338–15 349, 2020.
- [23] D. Ding, K. Zhang, T. Basar, and M. Jovanovic, "Natural policy gradient primal-dual method for constrained markov decision processes," *Advances in Neural Information Processing Systems*, vol. 33, pp. 8378–8390, 2020.
- [24] C. Tessler, D. J. Mankowitz, and S. Mannor, "Reward constrained policy optimization," in *International Conference on Learning Representations*, 2018.
- [25] R. J. Aumann, "Mixed and behavior strategies in infinite extensive games," in *Advances in Game Theory.(AM-52), Volume 52*. Princeton University Press, 2016, pp. 627–650.
- [26] A. Abate, M. Prandini, J. Lygeros, and S. Sastry, "Probabilistic reachability and safety for controlled discrete time stochastic hybrid systems," *Automatica*, vol. 44, pp. 2724–2734, 11 2008.
- [27] O. Hernández-Lerma and J. B. Lasserre, *Discrete-time Markov control processes: basic optimality criteria*. Springer Science & Business Media, 2012, vol. 30.
- [28] D. P. Bertsekas, *Dynamic Programming and Optimal Control*, 4th ed. Belmont, MA, USA: Athena Scientific, 2015, vol. I.
- [29] N. Kariotoglou, M. Kamgarpour, T. H. Summers, and J. Lygeros, "The linear programming approach to reach-avoid problems for markov decision processes," *Journal of Artificial Intelligence Research*, vol. 60, no. 1, p. 263–285, 2017.
- [30] S. P. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [31] E. J. Anderson and P. Nash, *Linear programming in infinite-dimensional spaces*. John Wiley & Sons, Ltd., Chichester, 1987, theory and applications, A Wiley-Interscience Publication.
- [32] A. S. Nowak, "Universally measurable strategies in zero-sum stochastic games," *The Annals of Probability*, vol. 13, no. 1, pp. 269–287, 1985.
- [33] L. D. Brown and R. Purves, "Measurable selections of extrema," *The annals of statistics*, pp. 902–912, 1973.
- [34] M. Ono, Y. Kuwata, and J. Balaram, "Mixed-strategy chance constrained optimal control," in *2013 American Control Conference*. IEEE, 2013, pp. 4666–4673.
- [35] E. A. Feinberg and A. Shwartz, "Constrained discounted dynamic programming," *Mathematics of Operations Research*, vol. 21, no. 4, pp. 922–945, 1996.
- [36] D. P. De Farias and B. Van Roy, "The linear programming approach to approximate dynamic programming," *Operations research*, vol. 51, no. 6, pp. 850–865, 2003.
- [37] N. Schmid and J. Lygeros, "Probabilistic reachability and invariance computation of stochastic systems using linear programming," *IFAC-PapersOnLine*, vol. 56, no. 2, pp. 11 229–11 234, 2023.
- [38] S. Summers and J. Lygeros, "Verification of discrete time stochastic hybrid systems: A stochastic reach-avoid decision problem," *Automatica*, vol. 46, no. 12, pp. 1951–1961, 2010.
- [39] J. W. Pitchford, E. A. Codling, and D. Psarra, "Uncertainty and sustainability in fisheries and the benefit of marine protected areas," *Ecological Modelling*, vol. 207, no. 2-4, pp. 286–292, 2007.