

# Robust Optimization for Hybrid MDPs with State-dependent Noise

Anonymous  
Paper ID #1907

## Abstract

Recent advances in solutions to Hybrid MDPs with discrete and continuous state and action spaces have significantly extended the class of MDPs for which exact solutions can be derived, albeit at the expense of a restricted transition noise model. In this paper, we work around limitations of previous solutions by adopting a robust optimization approach in which Nature is allowed to adversarially determine transition noise within pre-specified confidence intervals. This allows one to derive an optimal policy with an arbitrary (user-specified) level of success probability and significantly extends the class of transition noise models for which Hybrid MDPs can be solved. This work also significantly extends results for the related “chance-constrained” approach in stochastic hybrid control to accommodate state-dependent noise. We demonstrate our approach working on a variety of hybrid MDPs taken from AI planning, operations research, and control theory, noting that this is the first time optimal robust solutions have been automatically derived for such problems.

## 1 Introduction

Many real-world sequential decision-making problems are naturally modeled with both discrete and continuous (hybrid) state and action spaces. When state transitions are stochastic, these problems can be modeled as Hybrid Markov Decision Processes (HMDPs), which have been studied extensively in AI planning [Boyan and Littman, 2001; Feng *et al.*, 2004; Li and Littman, 2005; Kveton *et al.*, 2006; Marecki *et al.*, 2007; Meuleau *et al.*, 2009; Zamani *et al.*, 2012] as well as control theory [Henzinger *et al.*, 1997; Hu *et al.*, 2000; De Schutter *et al.*, 2009] and operations research [Puterman, 1994]. However, all previous solutions to hybrid MDPs either take an approximation approach or restrict stochastic noise on continuous transitions to be state-independent or discretized (i.e., requiring continuous transitions to be a finite mixture over deterministic transitions).

Unfortunately, each of these assumptions can be quite limiting in practice when strong *a priori* guarantees on performance are required in the presence of general forms of state-

dependent noise. For example, in a UAV NAVIGATION problem [Blackmore *et al.*, 2011], a human controller must be aware of all positions from which a UAV with a given amount of fuel reserves can return to its landing strip with high probability of success given known areas of (state-dependent) turbulence and weather events. In a SPACE TELESCOPE CONTROL problem [Löhr *et al.*, 2012], one must carefully manage inertial moments and rotational velocities as the telescope maneuvers between different angular orientations and zoom positions, where noise margins increase when the telescope is in unstable positions (extended zooms). And in a RESERVOIR CONTROL problem, one must manage reservoir levels to ensure a sufficient water supply for a population while avoiding overflow conditions subject to uncertainty over daily rainfall amounts. In all of these problems, there is no room for error: a UAV crash, a space telescope spinning uncontrollably, or a flooded reservoir can all cause substantial physical, monetary, and/or environmental damage. What is needed are robust solutions to these problems that are cost-optimal while guaranteed not to exceed a prespecified margin of error.

To achieve cost-optimal robust solutions we build on ideas used in the chance-constrained control literature [Schwarm and Nikolaou, 1999; Li *et al.*, 2002; Ono and Williams, 2008; Blackmore *et al.*, 2011] that maintain confidence intervals on (multivariate) noise distributions and ensure that all reachable states are within these noise margins. However, previous methods restrict either to linear systems with Gaussian uncertainty and state-independent noise or resort to approximation techniques. Furthermore, as these works are all inherently focused on control from a given initial state, they are unable to prove properties such as *robust controllability*, i.e., what states have a policy that can achieve a given cost with high certainty over some horizon?

In this work, we adopt a robust optimization receding horizon control approach in which Nature is allowed to adversarially determine transition noise w.r.t. constrained non-deterministic transitions in HMDPs. This permits us to find optimal robust solutions for a wide range of non-deterministic HMDPs and allows us to answer questions of *robust controllability* in very general state-dependent continuous noise settings. Altogether, this work significantly extends previous results in both the HMDP literature in AI and robust hybrid control literature and permits the solution of a new class of robust HMDP control problems.

## 2 Non-deterministic Hybrid MDPs

We first formally introduce the framework of Hybrid (discrete and continuous) Markov decision processes with non-deterministic continuous noise (ND-HMDPs) by extending the HMDP framework of [Zamani *et al.*, 2012]. The optimal solution for this model is then defined via robust dynamic programming.

### 2.1 Factored Representation

An HMDP is modelled using state variables  $(\vec{b}, \vec{x}) = (b_1, \dots, b_a, x_1, \dots, x_c)$  where each  $b_i \in \{0, 1\}$  ( $1 \leq i \leq a$ ) represents a discrete boolean variable and each  $x_j \in \mathbb{R}$  ( $1 \leq j \leq c$ ) is continuous. To model continuous uncertainty in ND-HMDPs we additionally define intermediate noise variables  $\vec{n} = n_1, \dots, n_e$  where each  $n_l \in \mathbb{R}$  ( $1 \leq l \leq e$ ). Both discrete and continuous actions are represented in the set  $A = \{a_1(\vec{y}_1), \dots, a_p(\vec{y}_p)\}$  where each action  $a(\vec{y}) \in A$  references a (possibly empty) vector of continuous parameters  $\vec{y} \in \mathbb{R}^{|\vec{y}|}$ ; we say an action is discrete if it has no continuous parameters ( $|\vec{y}| = 0$ ), otherwise it is continuous.

Given a current state  $(\vec{b}, \vec{x})$  and next state  $(\vec{b}', \vec{x}')$  and an executed action  $a(\vec{y})$  at the current state, a real-valued reward function  $R(\vec{b}, \vec{x}, \vec{b}', \vec{x}', a, \vec{y})$  specifies the immediate reward obtained at the current state. The probability of the next state  $(\vec{b}', \vec{x}')$  is defined by a joint state transition model  $P(\vec{b}', \vec{x}' | \vec{b}, \vec{x}, a, \vec{y}, \vec{n})$  which depends on the current state, action and noise. In a factored setting, we do not typically represent the transition distribution jointly but rather we factorize it into a dynamic Bayes net (DBN) as follows:

$$P(\vec{b}', \vec{x}' | \vec{b}, \vec{x}, a, \vec{y}, \vec{n}) = \prod_{i=1}^a P(b'_i | \vec{b}, \vec{x}, \vec{b}', \vec{x}', a, \vec{y}, \vec{n}) \prod_{j=1}^c P(x'_j | \vec{b}, \vec{x}, \vec{b}', \vec{x}', a, \vec{y}, \vec{n}) \quad (1)$$

Here we allow synchronic arcs under the condition that the DBN forms a proper directed acyclic graph (DAG). For binary variables  $b_i$  ( $1 \leq i \leq a$ ),  $P(b'_i | \vec{b}, \vec{x}, \vec{b}', \vec{x}', a, \vec{y}, \vec{n})$  are defined as general conditional probability functions (CPFs), which are not necessarily tabular since they may condition on inequalities over continuous variables. For continuous variables  $x_j$  ( $1 \leq j \leq c$ ), the CPFs  $P(x'_j | \vec{b}, \vec{x}, \vec{b}', \vec{x}', a, \vec{y}, \vec{n})$  are represented with *piecewise linear equations* (PLEs) that may have piecewise conditions which are arbitrary logical combinations of  $\vec{b}$ ,  $\vec{b}'$  and linear inequalities over  $\vec{x}$ ,  $\vec{x}'$ , and  $\vec{n}$ . Examples of PLEs will follow shortly.

In general, we assume that for each intermediate continuous noise variable  $n_l$  ( $1 \leq l \leq e$ ) a non-deterministic noise interval constraint function  $N(n_l | \vec{b}, \vec{x}, a, \vec{y})$  has been defined that represents a range covering  $\alpha$  of the probability mass for  $n_l$  and evaluates to  $-\infty$  for legal values of  $n_l$  and  $+\infty$  otherwise. The reason for the  $\pm\infty$  evaluation is simple: in a robust solution to HMDPs with non-deterministic noise constraints, Nature will attempt to adversarially minimize the reward the agent can achieve and hence we let  $N(n_l | \vec{b}, \vec{x}, a, \vec{y})$  take the value  $+\infty$  for illegal values of  $n_l$  to ensure Nature will never choose illegal assignments of  $n_l$  when minimizing.

As an intuitive example, if  $P(n_l | \vec{b}, \vec{x}, a, \vec{y}) = \mathcal{N}(n_l; \mu; \sigma^2)$  is a simple Normal distribution with mean  $\mu$  and variance  $\sigma^2$  and we let  $\alpha = 0.95$  then we know that that the 95% of the probability mass lies within  $\mu \pm 2\sigma$ , hence

$$N(n_l | \vec{b}, \vec{x}, a, \vec{y}) = \begin{cases} \mu - 2\sigma \leq n_l \leq \mu + 2\sigma : & -\infty \text{ (legal)} \\ \text{otherwise} : & +\infty \text{ (illegal)} \end{cases}$$

To make the ND-HMDP framework concrete, we now introduce a running example used throughout the paper:

**Example (RESERVOIR CONTROL).** *The problem of maintaining maximal reservoir levels subject to uncertain amounts of rainfall is an important problem in operations research (OR) literature [Mahootchi, 2009; Yeh, 1985]. In one variant of this problem, a reservoir operator must make a daily decision to drain some water from a reservoir or not subject to weather forecasts over some time horizon. Specifically in a seven day period, we assume that the weather forecast calls for a substantial amount of rain on the fourth day and chances of less rain on the others. The objective of the reservoir operator is to avoid underflow or overflow conditions while maximizing reservoir capacity.*

Formally, we assume a state consisting of continuous reservoir level  $l_1 \in \mathbb{R}$  and 3 boolean variables  $\vec{b}$  to encode a time period of eight days. We have two actions  $a \in A = \{\text{drain}, \text{no} - \text{drain}\}$ . The reward function  $R$  is used to prevent overflow and underflow by assigning  $-\infty$  penalty to water levels outside of lower reserve and upper capacity limits and a reward for the amount of water stored at the end of the time step. For both  $a \in A$  this is formally defined as:

$$R(l_1, l'_1, \vec{b}, \vec{b}', a) = \begin{cases} (200 \leq l_1 \leq 4500) \wedge (200 \leq l'_1 \leq 4500) : l'_1 \\ \text{otherwise} : & -\infty \end{cases}$$

For the transition function, we assume that on each time step  $\vec{b}' = \vec{b} + 1$  (not shown) and the reservoir level changes according to the amount of outflow (2000 units of water on a drain and 0 units on a no-drain action) plus a noisy (uncertain) amount of rain  $n$ :

$$P(l'_1 | l_1, n, d_i, d'_i, a = \text{drain}) = \delta(l'_1 - (n + l_1 - 2000))$$

$$P(l'_1 | l_1, n, d_i, d'_i, a = \text{no} - \text{drain}) = \delta(l'_1 - (n + l_1))$$

The use of the  $\delta[\cdot]$  function here ensures that the continuous CPF over  $l'$  integrates to 1, which is crucial for defining a proper probability distribution. While these PLEs are deterministic note that all continuous noise in this framework enters via the non-deterministic noise variables in ND-HMDPs. The noisy level of rainfall  $n$  is state-dependent and legal intervals are defined as follows:

$$N(n | \vec{b}, l_1) = \begin{cases} \vec{b} = 4 \wedge (1200 \leq n \leq 2000) : & -\infty \\ \vec{b} = 4 \wedge (0 \leq n \leq 400) : & -\infty \\ \text{otherwise} : & +\infty \end{cases}$$

In short, on day four ( $\vec{b} = 4$ ) the amount of rain is expected to be between 1200 and 2000 units, whereas on the other days it is expected to be between 0 and 400 units.

A policy  $\pi(\vec{b}, \vec{x})$  specifies the action  $a(\vec{y}) = \pi(\vec{b}, \vec{x})$  to take at state  $(\vec{b}, \vec{x})$ . In a robust solution to HMDPs with non-deterministic noise constraints, an optimal sequence of finite

horizon policies  $\Pi^* = (\pi^{*,1}, \dots, \pi^{*,H})$  is desired such that given the initial state  $(\vec{b}_0, \vec{x}_0)$  at  $h = 0$  and a discount factor  $\gamma$ ,  $0 \leq \gamma \leq 1$ , the expected sum of discounted rewards over horizon  $h \in H$  ( $H \geq 0$ ) is maximized subject to Nature's adversarial attempt to choose value minimizing assignments of the noise variables. The value function  $V$  w.r.t.  $\Pi^*$  in this case is defined via a recursive expectation

$$V^{\Pi^*,H}(\vec{b}, \vec{x}) = \min_{\vec{n}} \max \left( N(n_1 | \vec{b}, \vec{x}, \Pi^{*,H}), \dots, \max \left( N(n_e | \vec{b}, \vec{x}, \Pi^{*,H}), E_{\Pi^{*,H}} \left[ r^h + \gamma V^{\Pi^*,H-1}(\vec{b}', \vec{x}') \mid \vec{b}_0, \vec{x}_0 \right] \right) \dots \right)$$

where  $r^h$  is the reward obtained at horizon  $h$  following policy  $\Pi^*$  and using Nature's minimizing choice of  $\vec{n}$  at each  $h$ .

The effect of "max'ing" in each of the previously defined  $N(n_l | \vec{b}, \vec{x}, a, \vec{y})$  ( $1 \leq l \leq e$ ) with the value function is one of the major insights and contributions of this paper. We noted before that Nature will never choose an illegal value of  $n_l$  where  $N(n_l | \vec{b}, \vec{x}, a, \vec{y}) = +\infty$ , instead it will choose a legal value of  $n_l$  for which  $N(n_l | \vec{b}, \vec{x}, a, \vec{y}) = -\infty$  which when "max'ed" in with the value function effectively vanishes owing to the identity  $\max(v, -\infty) = v$  for all  $v > -\infty$ .

Finally, by leveraging the simple union bound, we can easily prove that a policy will achieve  $V^{\Pi^*,H}$  with at least  $1 - H(1 - \alpha)$  probability since the probability of encountering a noise value outside the confidence interval is only  $(1 - \alpha)$  at any time step. Hence for a success probability of at least  $\beta$ , one should choose  $\alpha = 1 - \frac{1-\beta}{H}$ , e.g.,  $\beta = 0.95$  success probability requires an  $\alpha = 0.99$  for  $H = 5$ .

## 2.2 Robust Dynamic Programming

We extend the value iteration dynamic programming algorithm [Bellman, 1957] and specifically the form used for HMDPs in [Zamani *et al.*, 2012] to a robust dynamic programming (RDP) algorithm for ND-HMDPs that may be considered a continuous action generalization of zero-sum alternating turn Markov games [Littman, 1994]. Initializing  $V^0(\vec{b}, \vec{x}) = 0$  the algorithm builds the  $h$ -stage-to-go value function  $V^h(\vec{b}, \vec{x})$ .

The quality  $Q_a^h(\vec{b}, \vec{x}, \vec{y}, \vec{n})$  of taking action  $a(\vec{y})$  in state  $(\vec{b}, \vec{x})$  with noise parameters  $\vec{n}$  and acting so as to obtain  $V^{h-1}(\vec{b}', \vec{x}')$  thereafter is defined as the following:

$$Q_a^h(\vec{b}, \vec{x}, \vec{y}, \vec{n}) = \max \left( N(n_1 | \vec{b}, \vec{x}, \Pi^{*,H}), \dots, \max \left( N(n_e | \vec{b}, \vec{x}, \Pi^{*,H}), \sum_{\vec{b}'} \prod_{i=1}^a P(b_i' | \vec{b}, \vec{x}, \vec{b}', \vec{x}', a, \vec{y}, \vec{n}) \prod_{j=1}^c P(x_j' | \vec{b}, \vec{x}, \vec{b}', \vec{x}', a, \vec{y}, \vec{n}) \left[ R(\vec{b}, \vec{x}, \vec{b}', \vec{x}', a, \vec{y}) + \gamma V^{h-1}(\vec{b}', \vec{x}') \right] \dots \right) \right) \quad (2)$$

Here the noise constraints  $N(\vec{n} | \vec{b}, \vec{x})$  are "max'ed" in with the value function to ensure Nature chooses a legal setting of  $n_l$ , effectively reducing each max to an identity operation.

Next, given  $Q_a^h(\vec{b}, \vec{x}, \vec{y}, \vec{n})$  as above for each  $a \in A$ , we can proceed to define the  $h$ -stage-to-go value function assuming that the agent attempts to maximize value subject to Nature's adversarial choice of value-minimizing noise:

$$V^h(\vec{b}, \vec{x}) = \max_{a \in A} \max_{\vec{y} \in \mathbb{R}^{|\vec{y}|}} \min_{\vec{n} \in \mathbb{R}^{|\vec{n}|}} \left\{ Q_a^h(\vec{b}, \vec{x}, \vec{y}, \vec{n}) \right\} \quad (3)$$

The optimal policy at horizon  $h$  can also be determined using the  $Q$ -function as below:

$$\pi^{*,h}(\vec{b}, \vec{x}) = \arg \max_{a \in A} \arg \max_{\vec{y} \in \mathbb{R}^{|\vec{y}|}} \min_{\vec{n} \in \mathbb{R}^{|\vec{n}|}} Q_a^h(\vec{b}, \vec{x}, \vec{y}, \vec{n}) \quad (4)$$

For finite-horizon HMDPs the optimal value function and policy are obtained up to horizon  $H$ . For infinite horizons where the optimal policy has finitely bounded value then value iteration terminates when two values are equal in subsequent horizons ( $V^h = V^{h-1}$ ). In this case  $V^\infty = V^h$  and  $\pi^{*,\infty} = \pi^{*,h}$ .

Up to this point we have only provided the abstract mathematical framework for ND-HMDPs and RDP. Fortunately though, we can leverage the continuous max (and analogously defined min) operations and symbolic DP approach of [Zamani *et al.*, 2012] in order to compute RDP via (2) and (3) exactly in closed-form. We discuss this next.

## 3 Robust Symbolic Dynamic Programming

In order to compute the equations above, we propose a *robust symbolic dynamic programming* (RSDP) approach building on the work of [Zamani *et al.*, 2012; Sanner *et al.*, 2011]. This requires a value iteration algorithm described in Algorithm 1 (VI) and the regression subroutine described in Algorithm 2. In what follows we show how the techniques of SDP can be extended to compute RDP exactly in closed-form as discussed in the last section.

In general we define *all* symbolic functions to be represented in *case* form [Boutillier *et al.*, 2001] for which a binary "cross-sum" operation can be defined as follows:

$$\left\{ \begin{array}{l} \phi_1 : f_1 \\ \phi_2 : f_2 \end{array} \right\} \oplus \left\{ \begin{array}{l} \psi_1 : g_1 \\ \psi_2 : g_2 \end{array} \right\} = \left\{ \begin{array}{l} \phi_1 \wedge \psi_1 : f_1 + g_1 \\ \phi_1 \wedge \psi_2 : f_1 + g_2 \\ \phi_2 \wedge \psi_1 : f_2 + g_1 \\ \phi_2 \wedge \psi_2 : f_2 + g_2 \end{array} \right\}$$

Here  $\phi_i$  and  $\psi_j$  are logical formulae defined over the state  $(\vec{b}, \vec{x})$  and can include arbitrary logical ( $\wedge, \vee, \neg$ ) combinations of boolean variables and *linear* inequalities ( $\geq, >, \leq, <$ ) over continuous variables – we call this *linear case form* (LCF). The  $f_i$  and  $g_j$  are restricted to be *linear* functions. Similarly operations such as  $\ominus$  and  $\otimes$  may be defined with operations applied to LCF functions yielded LCF results.

In addition to  $\ominus$  and  $\otimes$  another key binary operation on case statements the preserves the LCF property is *symbolic case maximization*:

$$\text{casemax} \left( \left\{ \begin{array}{l} \phi_1 : f_1 \\ \phi_2 : f_2 \end{array} \right\}, \left\{ \begin{array}{l} \psi_1 : g_1 \\ \psi_2 : g_2 \end{array} \right\} \right) = \left\{ \begin{array}{l} \phi_1 \wedge \psi_1 \wedge f_1 > g_1 : f_1 \\ \phi_1 \wedge \psi_1 \wedge f_1 \leq g_1 : g_1 \\ \phi_1 \wedge \psi_2 \wedge f_1 > g_2 : f_1 \\ \phi_1 \wedge \psi_2 \wedge f_1 \leq g_2 : g_2 \\ \vdots \end{array} \right\}$$

To demonstrate how VI symbolically implements RDP, we compute  $V^1$  for the RESERVOIR CONTROL example. For both actions, the function  $Q_a^1$  is computed in line 6 using Algorithm 2 with the following operations for action *no-drain*:

- Priming  $V$  which indicates a *symbolically substitution* of  $V^0 = V^0 \sigma = 0$  where  $\sigma = \{d_i \setminus d_i', l_j \setminus l_j'\}$ .
- Since the reward function contains the primed variable  $l_1'$ , line 4 is performed (and not line 15) where  $Q = R(l_1, l_1', d_i, d_i', n, a)$ .

**Algorithm 1:**  $\text{VI}(\text{CSA-MDP}, H) \rightarrow (V^h, \pi^{*,h})$ 


---

```

1 begin
2    $V^0 := 0, h := 0$ 
3   while  $h < H$  do
4      $h := h + 1$ 
5     foreach  $a(\vec{y}) \in A$  do
6        $Q_a^h(\vec{y}, \vec{n}) := \text{Regress}(V^{h-1}, a, \vec{y})$ 
7        $Q_a^h(\vec{y}) := \min_{\vec{n}} Q_a^h(\vec{y}, \vec{n})$  //Stochastic min
8        $Q_a^h := \max_{\vec{y}} Q_a^h(\vec{y})$  // Continuous max
9        $V^h := \text{casemax}_a Q_a^h$  // casemax all  $Q_a$ 
10       $\pi^{*,h} := \arg \max_{(a, \vec{y})} Q_a^h(\vec{y})$ 
11      if  $V^h = V^{h-1}$  then
12        break // Terminate if early convergence
13
14    return  $(V^h, \pi^{*,h})$ 
15 end

```

---

**Algorithm 2:**  $\text{Regress}(V, a, \vec{y}) \rightarrow Q$ 


---

```

1 begin
2    $Q = \text{Prime}(V)$  // All  $b_i \rightarrow b'_i$  and all  $x_i \rightarrow x'_i$ 
3   if  $v'$  in  $R$  then
4      $Q := R(\vec{b}, \vec{b}', \vec{x}, \vec{x}', a, \vec{y}) \oplus (\gamma \cdot Q)$ 
5
6   foreach  $v'$  in  $Q$  do
7     if  $v' = x'_j$  then
8       //Continuous marginal integration
9        $Q := \int Q \otimes P(x'_j | \vec{b}, \vec{b}', \vec{x}, \vec{x}', a, \vec{y}, \vec{n}) dx'_j$ 
10    if  $v' = b'_i$  then
11      //Discrete marginal summation
12       $Q := [Q \otimes P(b'_i | \vec{b}, \vec{b}', \vec{x}, \vec{x}', a, \vec{y}, \vec{n})] |_{b'_i=1}$ 
13       $\oplus [Q \otimes P(b'_i | \vec{b}, \vec{b}', \vec{x}, \vec{x}', a, \vec{y}, \vec{n})] |_{b'_i=0}$ 
14
15    if  $\neg (v' \text{ in } R)$  then
16       $Q := R(\vec{b}, \vec{b}', \vec{x}, \vec{x}', a, \vec{y}) \oplus (\gamma \cdot Q)$ 
17
18    foreach  $n_l$  in  $Q$  do
19      // Sequence of max-in for noise variables
20       $Q_a^h(\vec{y}, \vec{n}) := \text{casemax}_{n_l}(Q, N(n_l, b_i, x_j))$ 
21    return  $Q$ 
22 end

```

---

- For boolean variables, regression is performed using  $f|_{b=v}$  (restriction operator) which assigns the value  $v \in \{0, 1\}$  to any occurrence of  $b$  in  $f$  – not applicable to the example. For continuous variables line 9 follows the rules of integration w.r.t. a  $\delta$  function [Sanner *et al.*, 2011] which simply yields a symbolic substitution:

$$\int f(x'_j) \otimes \delta[x'_j - h(\vec{z})] dx'_j = f(x'_j) \{x'_j / h(\vec{z})\}$$

This results in the following  $Q$ -value for RESERVOIR

CONTROL :

$$\begin{cases} (200 \leq l_1 \leq 4500) \wedge (200 \leq (l_1 + n) \leq 4500) & : l_1 + n \\ \text{otherwise} & : -\infty \end{cases}$$

- Maximizing the result with each of the noise variables in defined line 20 using a sequence of symbolic maximizations. Each noise variable assigns  $-\infty$  for legal values inside the boundary range  $+\infty$  for illegal values defined by the noise model  $N(\vec{n}, \vec{b}, \vec{x})$ . The result is defined below:

$$\begin{cases} ((l_1 \wedge (l_1 + n)) \in \text{safe}) \wedge (n \in \text{legal}) & : l_1 + n \\ (l_1 \in \text{safe}) \wedge ((l_1 + n) \notin \text{safe}) \wedge (n \in \text{legal}) & : -\infty \\ (n \notin \text{legal}) & : +\infty \end{cases}$$

where *legal* noise value corresponds to  $[0, 400]$  or  $[1200, 2000]$  and *safe* water levels is  $[200, 4500]$ .

The regressed stochastic  $Q_a^h(\vec{y}, \vec{n})$  from Algorithm 2 is now minimized over the noise variables  $\vec{n}$  in line 7. Intuitively this continuous minimization will never choose  $+\infty$  as there is always some value smaller which insures that the transitioned model never chooses illegal values. Each partition  $i$  of this intermediate  $Q$  is considered for a continuous minimization separately with the final result a casemin (definition follows from casemax) on all the individual minimum results:  $\text{casemin}_i \min_n \phi_i(\vec{b}, \vec{x}, \vec{n}) f_i(\vec{b}, \vec{x}, \vec{n})$ . We demonstrate the steps of this algorithm for the third partition of the regressed  $Q$  defined as:

$$d_1 \wedge (200 \leq l_1 \leq 4500) \wedge (1200 \leq n \leq 2000) \wedge (200 \leq (l_1 + n) \leq 4500) : l_1 + n$$

For each partition the logical constraints are used to derive the (a) lower bound on  $n$  ( $LB = 1200, 200 - l_1$ ); (b) upper bound on  $n$  ( $UB = 2000, 4500 - l_1$ ) and (c) constraints independent of  $n$  ( $IND = d_1, 200 \leq l_1 \leq 4500$ ). In case of several bounds on  $n$  the maximum of all lower bounds and the minimum of all upper bounds is desired.:

$$LB = \begin{cases} l_1 < -1000 : 200 - l_1 \\ l_1 > -1000 : 1200 \end{cases}$$

$$UB = \begin{cases} l_1 > 2500 : 4500 - l_1 \\ l_1 < 2500 : 2000 \end{cases}$$

The minima points of upper and lower bounds are evaluated for the leaf value which equals to substituting the bounds instead of the noise variable  $n$  in the leaf function  $n + l_1$ :

$$Q = \begin{cases} l_1 \leq -1798 : 2000 + l_1 \\ l_1 \leq -1000 : 200 \\ l_1 \leq 3300 : 1200.48 + l_1 \\ l_1 \geq 3300 : 4500 \end{cases}$$

Natural constraints on bounds  $LB \leq UB$  and the *IND* constraints are also considered for the minimization on a single partition to obtain:

$$Q = \begin{cases} d_1 \wedge (2000 \leq l_1 \leq 3300) : 1200.48 + l_1 \\ \text{otherwise} : +\infty \end{cases}$$

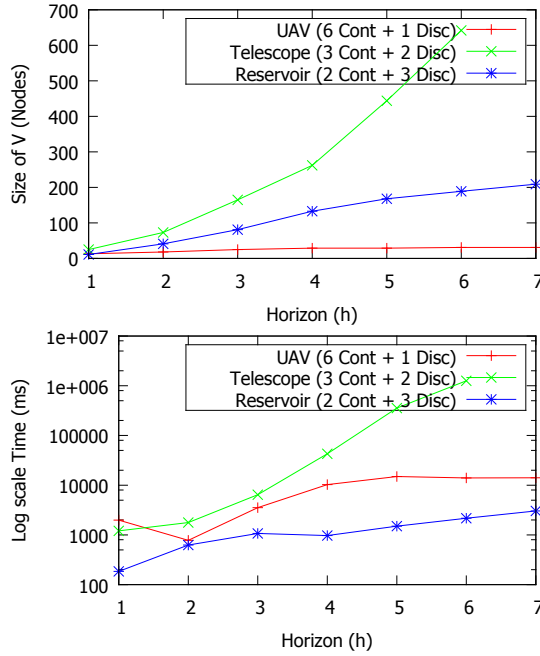


Figure 2: Space and elapsed time (between current and previous horizon) vs. horizon.

The final result of a continuous minimization is a casemin over all partitions which results in the following Q-value:

$$Q_{no-drain}^1 = \begin{cases} d_1 \wedge (200 \leq l_1 \leq 2506) : & 1200.48 + l_1 \\ \neg d_1 \wedge (200 \leq l_1 \leq 4098) : & l_1 \\ \text{otherwise} : & -\infty \end{cases}$$

The resulting Q-value with minimal noise is maximized over the continuous action parameter (not available in our example) in line 8. A discrete casemax on the set of discrete actions for all Q-functions defines the final V:

$$V^1 = \begin{cases} d_1 \wedge (2506 \leq l_1 \leq 4484) : & 1 \\ d_1 \wedge (200 \leq l_1 \leq 2506) : & 1200.48 + l_1 \\ \neg d_1 \wedge (4098 \leq l_1 \leq 4504) : & 1 \\ \neg d_1 \wedge (200 \leq l_1 \leq 4098) : & l_1 \\ \text{otherwise} : & -\infty \end{cases}$$

To implement the case statements efficiently with continuous variables, extended Algebraic Decision diagrams (XADDs) are used from [Sanner *et al.*, 2011] extended from ADDs [Bahar *et al.*, 1993].

In summary we remark that all operations including the continuous max and min operations preserve the LCF property, hence all operations for robust SDP can be performed exactly in closed-form — a first for receding horizon control with general forms of state-dependent continuous noise.

## 4 Empirical Results

We evaluated RH-MDP on the RESERVOIR CONTROL problem used as a running example, a UAV NAVIGATION problem and a SPACE TELESCOPE CONTROL problem — all

highly risk-sensitive and uncertain decision-making problems as described below.<sup>1</sup>

Figure 1 (left) shows the value function for the RESERVOIR CONTROL problem in horizon seven. We can observe that we can gain approximately 11000 units of reward if the reservoir begins with a water level approximately equal to 3000 at day zero. Otherwise, a lower initial starting state or the need to drain when the reservoir is near full lead to lower rewards for all states. Discontinuities in the value function occur at critical points where the policy changes over the time horizon.

**SPACE TELESCOPE CONTROL:** We have extended the problem of slewing a space telescope in order to look a new objective as given in [Löhr *et al.*, 2012]. This problem has six actions  $a_0, \dots, a_5$  that change the continuous angle  $k$  and angular rate  $v$ . The problem has one boolean state variable  $z$  for the telescope zoom state and one continuous noise variable. To model noise in this problem, we have only modified the transition function for the  $a_5$  action in the description from [Löhr *et al.*, 2012] (which could not handle noise) to add noise when  $v < 1 \frac{\text{deg}}{\text{seg}}$  and  $z = \text{false}$ :

$$\begin{aligned} k' &= (k + 40.55 * v) \\ v' &= (2/3v + n) \\ z' &= (\text{true}), \end{aligned}$$

We assume a noise in the transition function of the angular rate for  $a_5$  (which changes the zoom of the telescope and for which the dynamical model is only approximate) as follows:

$$n = \begin{cases} \neg(z) \wedge (n \leq 0.04 * v) \wedge (n \geq -0.04 * v) & : -\infty \\ \text{else} & : +\infty \end{cases}$$

which we note depends linearly on the angular velocity.

The reward for actions  $a_0, \dots, a_5$  is given by

$$R = \begin{cases} (z) \wedge (v \leq 0.02) \wedge (k \leq 1.683) \wedge (v \geq -0.02) \wedge (k \geq 1.283) & : 100 \\ \text{else} & : -\text{cost}(a) \end{cases}$$

where the  $\text{cost}(a)$  of action  $a_0$  is 0, 1 for actions  $a_i$   $i \in \{1, 2, 3, 4\}$  and 10 for action  $a_5$ . Figure 1 (right) shows the value function for the horizon four. We can see that there are relative few states that have a policy to achieve a goal (a reward of 0) with high certainty over this horizon.

Figure 1 (middle) shows the value function for the horizon four. We can observe that there are states with low angular rate ( $-0.04 \leq v \leq 0.04$  approximately) that have a policy to achieve a goal (a reward of 100) with high certainty over this horizon.

**UAV NAVIGATION:** In this problem a UAV needs to be able to plan trajectories that take the aircraft from its current location to a goal given constraints on time or fuel consumption and known areas of state-dependent turbulence (e.g., from localized weather events).

The state consist of UAVs continuous position  $x$  and  $y$ . In a given time step, the UAV may move a continuous distance  $ax \in [-40, 40]$  and  $ay \in [-40, 40]$ . The turbulence introduces a noise  $n_x$  and  $n_y$  respectively in the movement, given

<sup>1</sup>We note that all Java source code and a human/machine readable file format for all domains needed to reproduce the results in this paper can be found publicly online at *link suppressed to maintain anonymity*.

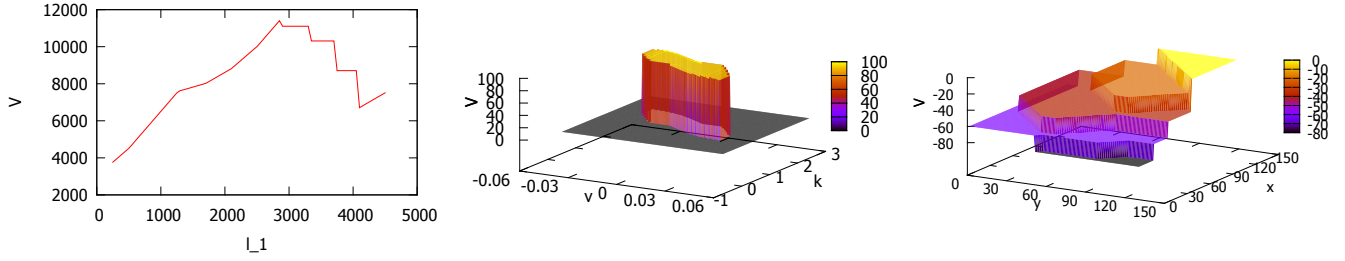


Figure 1: (left)  $V^7(l_1, d_1 = \text{false}, d_2 = \text{false}, d_3 = \text{false})$  RESERVOIR CONTROL problem; (middle)  $V^4(k, v, z = \text{true}, g = \text{false})$  SPACE TELESCOPE CONTROL problem. (right)  $V^4(x, y, l = \text{false})$  UAV NAVIGATION problem;

by:

$$n_x = \begin{cases} (y \geq 50 + x) \wedge (n_x \leq -20) \wedge (n_x \geq 20) & : -\infty \\ (y < 50 + x) \wedge (n_x \leq -5) \wedge (n_x \geq 5) & : -\infty \\ \text{else} & : +\infty \end{cases}$$

$$n_y = \begin{cases} (y \geq 50 + x) \wedge (n_y \leq -20) \wedge (n_y \geq 20) & : -\infty \\ (y < 50 + x) \wedge (n_y \leq -5) \wedge (n_y \geq 5) & : -\infty \\ \text{else} & : +\infty \end{cases}$$

The UAV goal is to achieve the region  $x + y > 200$ . It receives a reward penalty ( $-\infty$ ) for being in positions from which a UAV with a given amount of fuel reserves cannot return to its landing strip with high certainty. If the UAV is not in the goal position ( $\neg l$ ), the action reward is a cost of -20 fuel units for the given time period. We note that with six continuous variables in the regression (2 state, 2 action, 2 noise), this problem is relatively high-dimensional and could not be easily solved via discretization, which would also introduces error that we do not get in our optimal solution.

$$R = \begin{cases} (l) \wedge (x \leq 130) \wedge (y \leq 130) \wedge (x \geq 0) \wedge (y \geq 0) & : 0 \\ (\neg l) \wedge (x \leq 130) \wedge (y \leq 130) \wedge (x \geq 0) \wedge (y \geq 0) & : -20 \\ \text{else} & : -\infty \end{cases}$$

Figure 1 (right) shows the value function for the horizon four. We can see that there are relative few states that have a policy to achieve a goal (a reward of 0) with high certainty over this horizon.

Figure 2 shows the time and space for each of the solved problems. The UAV NAVIGATION problem has more continuous variables, however we can see that it is easier to solve than the SPACE TELESCOPE CONTROL, one possible reason is that this last problem has more actions with more complex forms of linearly state dependent noise.

## 5 Related Work

This work extends results in HMDP in AI [Boyan and Littman, 2001; Feng *et al.*, 2004; Li and Littman, 2005; Kveton *et al.*, 2006; Marecki *et al.*, 2007; Meuleau *et al.*, 2009; Zamani *et al.*, 2012] and hybrid system control literature [Henzinger *et al.*, 1997; Hu *et al.*, 2000; De Schutter *et al.*, 2009] to handled state-dependent noise.

In the hybrid control literature, a challenging topic is to solve the controllability problem that is NP hard [Blondel and Tsitsiklis, 1999]. A hybrid system is called hybrid controllable if, for any pair of valid states, there exists at

least one permitted control sequence (correct control-laws) between them [Tittus and Egardt, 1998; Yang and Blanke, 2007]. Another challenging topic for stochastic hybrid systems, a class of hybrid systems that allows uncertainty, is tried to maximize the probability that the execution will remain in safe states as long as possible [Hu *et al.*, 2000]. This work is related with both topics, however we want to answer a slightly different question, called the robust controllability problem: what states have a policy to achieve a goal (that can be modeled as a reward or cost function) with high certainty over some horizon? To the authors knowledge, in the control area there are few results to answer a similar question except in the chance-constrained predictive stochastic sub-area, that finds the optimal sequence of control inputs subject to the constraint that the probability of failure must be below a user-specified threshold [Blackmore *et al.*, 2011]. However all the previous work in this sub-area is focused on linear systems subject to Gaussian uncertainty and state-independence noise [Schwarm and Nikolaou, 1999; Li *et al.*, 2002; Ono and Williams, 2008; Blackmore *et al.*, 2011] or resort to approximation techniques [Blackmore *et al.*, 2010]. We remark that our approach is not approximated and can optimally solve problems with state-dependent noise in a receding horizon control framework that answers the robust controllability question.

## 6 Concluding Remarks

This work has combined symbolic techniques and data structures from the HMDP literature in AI with techniques from chance-constrained control theory to provide optimal robust solutions to a range of problems with general continuous transitions and state-dependent noise for which no general exact closed-form solutions previously existed. Using these techniques we were able to find optimal policies and answer questions of robust controllability for a variety of highly risk-sensitive applications from AI planning, control theory, and operations research such as UAV NAVIGATION, SPACE TELESCOPE CONTROL, and RESERVOIR CONTROL. Among many potential avenues for future work, combining this receding horizon control approach with focused search techniques as in HAO\* [Meuleau *et al.*, 2009] should preserve our strong robust optimality guarantees while substantially increasing the scalability of our approach in exchange for restricting solution optimality to a known set of initial states.

## References

- [Bahar *et al.*, 1993] R. Iris Bahar, Erica Frohm, Charles Gaona, Gary Hachtel, Enrico Macii, Abelardo Pardo, and Fabio Somenzi. Algebraic Decision Diagrams and their applications. In *IEEE /ACM International Conference on CAD*, 1993.
- [Bellman, 1957] Richard E. Bellman. *Dynamic Programming*. Princeton University Press, Princeton, NJ, 1957.
- [Blackmore *et al.*, 2010] L. Blackmore, M. Ono, A. Bekasov, and B.C. Williams. A probabilistic particle-control approximation of chance-constrained stochastic predictive control. *Robotics, IEEE Transactions on*, 26(3):502–517, 2010.
- [Blackmore *et al.*, 2011] Lars Blackmore, Masahiro Ono, and Brian C. Williams. Chance-constrained optimal path planning with obstacles. *IEEE Transactions on Robotics*, 27(6):1080–1094, 2011.
- [Blondel and Tsitsiklis, 1999] Vincent D. Blondel and John N. Tsitsiklis. Complexity of stability and controllability of elementary hybrid systems. *Automatica*, 35(3):479–489, 1999.
- [Boutilier *et al.*, 2001] Craig Boutilier, Ray Reiter, and Bob Price. Symbolic dynamic programming for first-order MDPs. In *IJCAI-01*, pages 690–697, Seattle, 2001.
- [Boyan and Littman, 2001] Justin Boyan and Michael Littman. Exact solutions to time-dependent MDPs. In *Advances in Neural Information Processing Systems NIPS-00*, pages 1026–1032, 2001.
- [De Schutter *et al.*, 2009] B. De Schutter, W.P.M.H. Heemels, J. Lunze, and C. Prieur. Survey of modeling, analysis, and control of hybrid systems. In J. Lunze and F. Lamnabhi-Lagarigue, editors, *Handbook of Hybrid Systems Control – Theory, Tools, Applications*, chapter 2, pages 31–55. Cambridge University Press, Cambridge, UK, 2009.
- [Feng *et al.*, 2004] Zhengzhu Feng, Richard Dearden, Nicolas Meuleau, and Richard Washington. Dynamic programming for structured continuous markov decision problems. In *Uncertainty in Artificial Intelligence (UAI-04)*, pages 154–161, 2004.
- [Henzinger *et al.*, 1997] Thomas A. Henzinger, Pei H. Ho, and Howard W. Toi. HYTECH: A Model Checker for Hybrid Systems. *International Journal on Software Tools for Technology Transfer*, 1(1-2):110–122, 1997.
- [Hu *et al.*, 2000] J. Hu, John Lygeros, and S. Sastry. Towards a theory of stochastic hybrid systems. *Lecture Notes in Computer Science LNCS*, 1790:160–173, 2000.
- [Kveton *et al.*, 2006] Branislav Kveton, Milos Hauskrecht, and Carlos Guestrin. Solving factored mdps with hybrid state and action variables. *Journal Artificial Intelligence Research (JAIR)*, 27:153–201, 2006.
- [Li and Littman, 2005] Lihong Li and Michael L. Littman. Lazy approximation for solving continuous finite-horizon mdps. In *National Conference on Artificial Intelligence AAAI-05*, pages 1175–1180, 2005.
- [Li *et al.*, 2002] Pu Li, Moritz Wendt, and Günter Wozny. Brief a probabilistically constrained model predictive controller. *Automatica*, 38(7):1171–1176, 2002.
- [Littman, 1994] Michael L. Littman. Markov games as a framework for multi-agent reinforcement learning. In *ICML*, pages 157–163, 1994.
- [Löhr *et al.*, 2012] Johannes Löhr, Patrick Eyerich, Thomas Keller, and Bernhard Nebel. A planning based framework for controlling hybrid systems. In *Proceedings of the Twenty-Second International Conference on Automated Planning and Scheduling*, 2012.
- [Mahootchi, 2009] Masoud Mahootchi. *Storage System Management Using Reinforcement Learning Techniques and Nonlinear Models*. PhD thesis, University of Waterloo, Canada, 2009.
- [Marecki *et al.*, 2007] Janusz Marecki, Sven Koenig, and Milind Tambe. A fast analytical algorithm for solving markov decision processes with real-valued resources. In *International Conference on Uncertainty in Artificial Intelligence IJCAI*, pages 2536–2541, 2007.
- [Meuleau *et al.*, 2009] Nicolas Meuleau, Emmanuel Benazera, Ronen I. Brafman, Eric A. Hansen, and Mausam. A heuristic search approach to planning with continuous resources in stochastic domains. *Journal Artificial Intelligence Research (JAIR)*, 34:27–59, 2009.
- [Ono and Williams, 2008] Masahiro Ono and Brian C. Williams. An efficient motion planning algorithm for stochastic dynamic systems with constraints on probability of failure. In *AAAI*, pages 1376–1382, 2008.
- [Puterman, 1994] Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley, New York, 1994.
- [Sanner *et al.*, 2011] Scott Sanner, Karina Valdivia Delgado, and Leliane Nunes de Barros. Symbolic dynamic programming for discrete and continuous state mdps. In *Proceedings of the 27th Conference on Uncertainty in AI (UAI-2011)*, Barcelona, 2011.
- [Schwarm and Nikolaou, 1999] Alexander T. Schwarm and Michael Nikolaou. Chance-constrained model predictive control. *AIChE Journal*, 45(8):1743–1752, 1999.
- [Tittus and Egardt, 1998] M. Tittus and B. Egardt. Control design for integrator hybrid systems. *Automatic Control, IEEE Transactions on*, 43(4):491–500, apr 1998.
- [Yang and Blanke, 2007] Zhenyu Yang and Mogens Blanke. A unified approach to controllability analysis for hybrid control systems. *Nonlinear Analysis: Hybrid Systems*, 1(2):212–222, 2007.
- [Yeh, 1985] William G Yeh. Reservoir management and operations models: A state-of-the-art review. *Water Resources research*, 21,12:17971818, 1985.
- [Zamani *et al.*, 2012] Z. Zamani, S. Sanner, and C. Fang. Symbolic dynamic programming for continuous state and action mdps. In *In Proceedings of the 26th AAAI Conference on Artificial Intelligence (AAAI-12)*, Toronto, Canada, 2012.