
Symbolic Dynamic Programming for Continuous State and Observation POMDPs

Zahra Zamani
ANU & NICTA
Canberra, Australia
zahra.zamani@anu.edu.au

Scott Sanner
NICTA & ANU
Canberra, Australia
scott.sanner@nicta.com.au

Pascal Poupart
U. of Waterloo
Waterloo, Canada
ppoupart@uwaterloo.ca

Kristian Kersting
Fraunhofer IAIS & U. of Bonn
Bonn, Germany
kristian.kersting@iais.fraunhofer.de

Abstract

Point-based value iteration (PBVI) methods have proven extremely effective for finding (approximately) optimal dynamic programming solutions to partially-observable Markov decision processes (POMDPs) when a set of initial belief states is known. However, no PBVI work has provided *exact point-based backups for both continuous state and observation spaces*, which we tackle in this paper. Our key insight is that while there may be an infinite number of observations, there are only a finite number of continuous observation partitionings that are relevant for optimal decision-making when a finite, fixed set of reachable belief states is considered. To this end, we make two important contributions: (1) we show how previous exact symbolic dynamic programming solutions for continuous state MDPs can be generalized to *continuous state POMDPs with discrete observations*, and (2) we show how recently developed symbolic integration methods allow this solution to be extended to PBVI for *continuous state and observation POMDPs* with potentially correlated, multivariate continuous observation spaces.

1 Introduction

Partially-observable Markov decision processes (POMDPs) are a powerful modeling formalism for real-world sequential decision-making problems [3]. In recent years, point-based value iteration methods (PBVI) [5, 10, 11, 7] have proved extremely successful at scaling (approximately) optimal POMDP solutions to large state spaces when a set of initial belief states is known.

While PBVI has been extended to both continuous state and continuous observation spaces, no prior work has tackled both jointly without sampling. [6] provides exact point-based backups for continuous state and discrete observation problems (with approximate sample-based extensions to continuous actions and observations), while [2] provides exact point-based backups (PBBs) for discrete state and continuous observation problems (where multivariate observations must be conditionally independent). While restricted to discrete states, [2] provides an important insight that we exploit in this work: *only a finite number of partitionings of the observation space are required to distinguish between the optimal conditional policy over a finite set of belief states*.

We propose two major contributions: First, we extend symbolic dynamic programming for continuous state MDPs [9] to POMDPs with discrete observations, *arbitrary* continuous reward and transitions with discrete noise (i.e., a finite mixture of deterministic transitions). Second, we extend this symbolic dynamic programming algorithm to PBVI and the case of continuous observations

(while restricting transition dynamics to be piecewise linear with discrete noise, rewards to be piecewise constant, and observation probabilities and beliefs to be uniform) by building on [2] to *derive* relevant observation partitions for potentially correlated, multivariate continuous observations.

2 DC-POMDP Model

A discrete and continuous partially observable MDP (DC-POMDP) is a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T}, \mathcal{R}, \mathcal{Z}, \gamma, h \rangle$. States \mathcal{S} are given by vector $\mathbf{dx}_s = (\mathbf{d}_s, \mathbf{x}_s) = (d_{s_1}, \dots, d_{s_n}, x_{s_1}, \dots, x_{s_m})$ where each $d_{s_i} \in \{0, 1\}$ ($1 \leq i \leq n$) is boolean and each $x_{s_j} \in \mathbb{R}$ ($1 \leq j \leq m$) is continuous. We assume a finite, discrete action space $\mathcal{A} = \{a_1, \dots, a_r\}$. Observations \mathcal{O} are given by the vector $\mathbf{dx}_o = (\mathbf{d}_o, \mathbf{x}_o) = (d_{o_1}, \dots, d_{o_p}, x_{o_1}, \dots, x_{o_q})$ where each $d_{o_i} \in \{0, 1\}$ ($1 \leq i \leq p$) is boolean and each $x_{o_j} \in \mathbb{R}$ ($1 \leq j \leq q$) is continuous.

Three functions are required for modeling DC-POMDPs: (1) $\mathcal{T} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ a Markovian transition model defined as the probability of the next state given the action and previous state; (2) $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ a reward function which returns the immediate reward of taking an action in some state; and (3) an observation function defined as $\mathcal{Z} : \mathcal{S} \times \mathcal{A} \times \mathcal{O} \rightarrow [0, 1]$ which gives the probability of an observation given the outcome of a state after executing an action. A discount factor γ , $0 \leq \gamma \leq 1$ is used to discount rewards t time steps into the future by γ^t .

We use a dynamic Bayes net (DBN)¹ to compactly represent the transition model \mathcal{T} over the factored state variables and we use a two-layer Bayes net to represent the observation model \mathcal{Z} :

$$\mathcal{T} : p(\mathbf{dx}'_s | \mathbf{dx}_s, a) = \prod_{i=1}^n p(d'_{s_i} | \mathbf{dx}_s, a) \prod_{j=1}^m p(x'_{s_j} | \mathbf{dx}_s, \mathbf{d}'_s, a). \quad (1)$$

$$\mathcal{Z} : p(\mathbf{dx}_o | \mathbf{dx}'_s, a) = \prod_{i=1}^p p(d_{o_i} | \mathbf{dx}'_s, a) \prod_{j=1}^q p(x_{o_j} | \mathbf{dx}'_s, a). \quad (2)$$

Probabilities over *discrete* variables $p(d'_{s_i} | \mathbf{dx}_s, a)$ and $p(d_{o_i} | \mathbf{dx}'_s, a)$ may condition on both discrete variables and (nonlinear) inequalities of continuous variables; this is further restricted to linear inequalities in the case of continuous observations. Transitions over *continuous* variables $p(x'_{s_j} | \mathbf{dx}_s, \mathbf{d}'_s, a)$ must be deterministic (but arbitrary nonlinear) piecewise functions; in the case of continuous observations they are further restricted to be piecewise linear; this permits discrete noise in the continuous transitions since they may condition on stochastically sampled discrete next-state variables \mathbf{d}'_s . Observation probabilities over continuous variables $p(x_{o_j} | \mathbf{dx}'_s, a)$ only occur in the case of continuous observation and are required to be piecewise constant (a mixture of uniform distributions); the same restriction holds for belief state representations. The reward $R(\mathbf{d}, \mathbf{x}, a)$ may be an arbitrary (nonlinear) piecewise function in the case of deterministic observations and a piecewise constant function in the case of continuous observations. We now provide concrete examples.

Example (Power Plant) [1] *The steam generation system of a power plant evaporates feed-water under restricted pressure and temperature conditions to turn a steam turbine. A reward is obtained when electricity is generated from the turbine and the steam pressure and temperature are within safe ranges. Mixing water and steam makes the respective pressure and temperature observations $p_o \in \mathbb{R}$ and $t_o \in \mathbb{R}$ on the underlying state $p_s \in \mathbb{R}$ and $t_s \in \mathbb{R}$ highly uncertain. Actions $A = \{\text{open}, \text{close}\}$ control temperature and pressure by means of a pressure valve.*

We initially present two DC-POMDP variants labeled **1D-Power Plant** using a single temperature state variable t_s . The transition and reward are common to both — temperature increments (decrements) with a closed (opened) valve, a large negative reward is given for a closed valve with t_s exceeding critical threshold 15, and positive reward is given for a safe, electricity-producing state:

$$p(t'_s | t_s, a) = \delta \left[t'_s - \begin{cases} (a = \text{open}) & : t_s - 5 \\ (a = \text{close}) & : t_s + 7 \end{cases} \right] \quad R(t_s, a) = \begin{cases} (a = \text{open}) & : -1 \\ (a = \text{close}) \wedge (t_s > 15) & : -1000 \\ (a = \text{close}) \wedge \neg(t_s > 15) & : 100 \end{cases} \quad (3)$$

Next we introduce the **Discrete Obs. 1D-Power Plant** variant where we define an *observation space* with a single discrete binary variable $o \in \mathcal{O} = \{\text{high}, \text{low}\}$:

¹We disallow general synchronic arcs for simplicity of exposition but note their inclusion only places restrictions on the variable elimination ordering used during the dynamic programming backup operation.

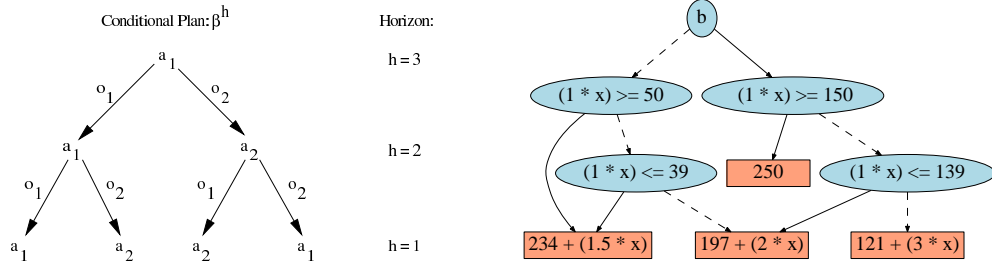


Figure 1: (left) Example conditional plan β^h for discrete observations; (right) example α -function for β^h over state $b \in \{0, 1\}$, $x \in \mathbb{R}$ in decision diagram form: the *true* (1) branch is solid, the *false* (0) branch is dashed.

$$p(o = \text{high} | t'_s, a = \text{open}) = \begin{cases} t'_s \leq 15 & : 0.9 \\ t'_s > 15 & : 0.1 \end{cases} \quad p(o = \text{high} | t'_s, a = \text{close}) = \begin{cases} t'_s \leq 15 & : 0.7 \\ t'_s > 15 & : 0.3 \end{cases} \quad (4)$$

Finally we introduce the **Cont. Obs. 1D-Power Plant** variant where we define an *observation space* with a single continuous variable t_o uniformly distributed on an interval of 10 units centered at t'_s .

$$p(t_o | t'_s, a = \text{open}) = U(t_o; t'_s - 5, t'_s + 5) = \begin{cases} (t_o > t'_s - 5) \wedge (t_o < t'_s + 5) & : 0.1 \\ (t_o \leq t'_s - 5) \vee (t_o \geq t'_s + 5) & : 0 \end{cases} \quad (5)$$

While simple, we note no prior method could perform exact point-based backups for either problem.

3 Value Iteration for DC-POMDPs

In a DC-POMDP, the agent does not directly observe the states and thus must maintain a belief state $b(\mathbf{dx}_s) = p(\mathbf{dx}_s)$. For a given belief state $\mathbf{b} = b(\mathbf{dx}_s)$, a POMDP policy π can be represented by a tree corresponding to a conditional plan β . An h -step conditional plan β^h can be defined recursively in terms of $(h - 1)$ -step conditional plans as shown in Fig. 1 (left). Our goal is to find a policy π that maximizes the value function, defined as the sum of expected discounted rewards over horizon h starting from initial belief state \mathbf{b} :

$$V_\pi^h(\mathbf{b}) = E_\pi \left[\sum_{t=0}^h \gamma^t \cdot r_t \mid \mathbf{b}_0 = \mathbf{b} \right] \quad (6)$$

where r_t is the reward obtained at time t and \mathbf{b}_0 is the belief state at $t = 0$. For finite h and belief state \mathbf{b} , the optimal policy π is given by an h -step conditional plan β^h . For $h = \infty$, the optimal discounted ($\gamma < 1$) value can be approximated arbitrarily closely by a sufficiently large h [3].

Even when the state is continuous (but the actions and observations are discrete), the optimal POMDP value function for finite horizon h is a piecewise linear and convex function of the belief state \mathbf{b} [6], hence V^h is given by a maximum over a finite set of “ α -functions” α_i^h :

$$V^h(\mathbf{b}) = \max_{\alpha_i^h \in \Gamma^h} \langle \alpha_i^h, \mathbf{b} \rangle = \max_{\alpha_i^h \in \Gamma^h} \int_{\mathbf{x}_s} \sum_{\mathbf{d}_s} \alpha_i^h(\mathbf{dx}_s) \cdot \mathbf{b}(\mathbf{dx}_s) d\mathbf{x}_s \quad (7)$$

Later on when we tackle continuous state *and* observations, we note that we will dynamically derive an optimal, finite partitioning of the observation space for a given belief state and hence reduce the continuous observation problem back to a discrete observation problem at every horizon.

The Γ^h in this optimal h -stage-to-go value function can be computed via Monahan’s dynamic programming approach to *value iteration* (VI) [4]. Initializing $\alpha_1^0 = \mathbf{0}$, $\Gamma^0 = \{\alpha_1^0\}$, and assuming discrete observations $o \in \mathcal{O}^h$, Γ^h is obtained from Γ^{h-1} as follows:²

$$g_{a,o,j}^h(\mathbf{dx}_s) = \int_{\mathbf{x}_{s'}} \sum_{\mathbf{d}_{s'}} p(o | \mathbf{dx}_{s'}, a) p(\mathbf{dx}_{s'} | \mathbf{dx}_s, a) \alpha_j^{h-1}(\mathbf{dx}_{s'}) d\mathbf{x}_{s'}; \quad \forall \alpha_j^{h-1} \in \Gamma^{h-1} \quad (8)$$

$$\Gamma_a^h = R(\mathbf{dx}_s, a) + \gamma \boxplus_{o \in \mathcal{O}} \left\{ g_{a,o,j}^h(\mathbf{dx}_s) \right\}_j \quad (9)$$

$$\Gamma^h = \bigcup_a \Gamma_a^h \quad (10)$$

²The \boxplus of sets is defined as $\boxplus_{j \in \{1, \dots, n\}} S_j = S_1 \boxplus \dots \boxplus S_n$ where the pairwise cross-sum $P \boxplus Q = \{\mathbf{p} + \mathbf{q} \mid \mathbf{p} \in P, \mathbf{q} \in Q\}$.

Algorithm 1: $\text{PBVI}(\text{DC-POMDP}, H, B = \{\mathbf{b}_i\}) \rightarrow \langle V^h \rangle$

```

1 begin
2    $V^0 := 0, h := 0, \Gamma_{PBVI}^0 = \{\alpha_1^0\}$ 
3   while  $h < H$  do
4      $h := h + 1, \Gamma^h := \emptyset, \Gamma_{PBVI}^h := \emptyset$ 
5     foreach  $\mathbf{b}_i \in B$  do
6       foreach  $a \in A$  do
7          $\Gamma_a^h := \emptyset$ 
8         if (continuous observations:  $q > 0$ ) then
9           // Derive relevant observation partitions  $\mathcal{O}_i^h$  for belief  $\mathbf{b}_i$ 
10           $\langle \mathcal{O}_i^h, p(\mathcal{O}_i^h | \mathbf{dx}_s', a) \rangle := \text{GenRelObs}(\Gamma_{PBVI}^{h-1}, a, \mathbf{b}_i)$ 
11        else
12          // Discrete observations and model already known
13           $\mathcal{O}_i^h := \{\mathbf{d}_o\}; p(\mathcal{O}_i^h | \mathbf{dx}_s', a) := \text{see Eq (2)}$ 
14          foreach  $o \in \mathcal{O}_i^h$  do
15            foreach  $\alpha_j^{h-1} \in \Gamma_{PBVI}^{h-1}$  do
16               $\alpha_j^{h-1} := \text{Prime}(\alpha_j^{h-1})$  //  $\forall d_i: d_i \rightarrow d_i'$  and  $\forall x_i: x_i \rightarrow x_i'$ 
17               $g_{a,o,j}^h := \text{see Eq (8)}$ 
18             $\Gamma_a^h := \text{see Eq (9)}$ 
19           $\Gamma^h := \Gamma^h \cup \Gamma_a^h$ 
20        foreach  $\mathbf{b}_i \in B$  do
21           $\alpha_{\mathbf{b}_i}^h := \arg \max_{\alpha_j \in \Gamma^h} \alpha_j \cdot \mathbf{b}_i$ 
22           $\Gamma_{PBVI}^h := \Gamma_{PBVI}^h \cup \alpha_{\mathbf{b}_i}^h$ 
23        if  $\Gamma_{PBVI}^h = \Gamma_{PBVI}^{h-1}$  then
24          break // Terminate if early convergence
25      return  $\Gamma_{PBVI}$ 
26 end

```

Point-based value iteration (PBVI) computes the value function only for a set of belief states $\{\mathbf{b}_i\}$ where $\mathbf{b}_i := p(\mathbf{dx}_s)$. The idea is straightforward and the main modification needed to Monahan's VI approach in Algorithm 1 (following [6]) is the loop from lines 22–24 where only α -vectors optimal at some belief state are retained for subsequent iterations. In the case of continuous observation variables ($q > 0$), we will need to derive a relevant set of observations on line 10, a key contribution of this work as described in Section 4.3. Whereas PBVI is optimal if all reachable belief states within horizon H are enumerated in B , in the DC-POMDP setting, the generation of continuous observations will most often lead to an infinite number of reachable belief states, even with finite horizon. Nonetheless, PBVI has been quite successful in practice without exhaustive enumeration of all reachable beliefs [5, 10, 11, 7], which motivates our use of PBVI in this work.

4 Symbolic Dynamic Programming

In this section we take a symbolic dynamic programming (SDP) approach to implementing VI and PBVI as defined in the last section. To do this, we need only show that all required operations can be computed efficiently and in closed-form, which we do next, building on SDP for MDPs [9].

4.1 Case Representation and Extended ADDs

The previous **Power Plant** examples represented all functions in case form, generally defined as

$$f = \begin{cases} \phi_1 : & f_1 \\ \vdots & \vdots \\ \phi_k : & f_k \end{cases}$$

and this is the form we use to represent all functions in a DC-POMDP. The ϕ_i are disjoint logical formulae defined over \mathbf{dx}_s and/or \mathbf{dx}_o with logical (\wedge, \vee, \neg) combinations of boolean variables and inequalities ($\geq, >, \leq, <$) over continuous variables. For discrete observation DC-POMDPs, the f_i and inequalities may use any function (e.g., $\sin(x_1) > \log(x_2) \cdot x_3$); for continuous observations, they are restricted to linear inequalities and linear or piecewise constant f_i as described in Section 2.

For *unary operations* such as scalar multiplication $c \cdot f$ (for some constant $c \in \mathbb{R}$) or negation $-f$ on case statements is simply to apply the operation on each case partition f_i ($1 \leq i \leq k$). A *binary operation* on two case statements, takes the cross-product of the logical partitions of each case statement and performs the corresponding operation on the resulting paired partitions. The cross-sum \oplus of two cases is defined as the following:

$$\begin{cases} \phi_1 : f_1 \\ \phi_2 : f_2 \end{cases} \oplus \begin{cases} \psi_1 : g_1 \\ \psi_2 : g_2 \end{cases} = \begin{cases} \phi_1 \wedge \psi_1 : f_1 + g_1 \\ \phi_1 \wedge \psi_2 : f_1 + g_2 \\ \phi_2 \wedge \psi_1 : f_2 + g_1 \\ \phi_2 \wedge \psi_2 : f_2 + g_2 \end{cases}$$

Likewise \ominus and \otimes are defined by subtracting or multiplying partition values. Inconsistent partitions can be discarded when they are irrelevant to the function value. A *symbolic case maximization* is defined as below:

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

The following SDP operations on case statements require more detail than can be provided here, hence we refer the reader to the relevant literature: *Restriction* $f|_\phi$: Takes a function f to restrict only in cases that satisfy some formula ϕ as defined in [9]. *Substitution* $f\sigma$: Takes a set σ of variables and their substitutions (which may be case statements), and carries out all variable substitutions in sequence [9]. *Integration* $\int_{x_1} f dx_1$: There are two forms: If x_1 is involved in a δ -function (cf. the transition in Eq (3)) then the integral is equivalent to a symbolic substitution and can be applied to any case statement [9]. Otherwise, if f is restricted to linear constraints and constant values, then the approach of [8] can be applied to yield a linearly constrained piecewise linear result.

The data structure of the *extended algebraic decision diagram* (XADD) [9] is used to support case statements and the required operations. Figure 1 (right) is an example of an XADD representation.

4.2 VI for DC State and Discrete Observations

For DC-POMDPs with only discrete observations $o \in \mathcal{O}$ and observation function $p(o|\mathbf{dx}'_s, a)$ (e.g., Eq (4)), we introduce a symbolic version of Monahan’s VI algorithm. In brief, we note that all VI operations needed in Section 3 apply *directly* to DC-POMDPs, e.g., we can rewrite Eq (8):

$$g_{a,o,j}^h(\mathbf{dx}_s) = \int_{\mathbf{x}_{s'} \mathbf{d}_{s'}} \left[p(o|\mathbf{dx}'_s, a) \otimes \left(\bigotimes_{i=1}^n p(d'_{s_i}|\mathbf{dx}_s, a) \right) \otimes \left(\bigotimes_{j=1}^m p(x'_{s_j}|\mathbf{dx}_s, \mathbf{d}'_s, a) \right) \otimes \alpha_j^{h-1}(\mathbf{dx}'_s) \right] d\mathbf{x}_{s'} \quad (11)$$

Crucially we note since the continuous transition cpfs $p(x'_{s_j}|\mathbf{dx}_s, \mathbf{d}'_s, a)$ are deterministic and hence defined with Dirac δ ’s (e.g., Eq 3) as described in Section 2, the integral $\int_{\mathbf{x}_{s'}}$ can always be computed in closed case form as discussed in Section 4.1. In short, nothing additional is required for PBVI on DC-POMDPs in this case — the key insight is simply that α -functions are now represented by case statements and can “grow” with the horizon as they partition the state space more and more finely.

4.3 PBVI for DC State and DC Observations

In general, it would be impossible to apply standard VI to DC-POMDPs with continuous observations since the number of observations is infinite. However, building on ideas in [2], in the case of PBVI, it is possible to *derive* a finite set of continuous observation partitions that permit exact point-based backups *at a belief point*. This additional operation (GenRelObs) appears on line 10 of PBVI in Algorithm 1 in the case of continuous observations and is formally defined in Algorithm 2.

To demonstrate the generation of relevant continuous observation partitions, we use the second iteration of the **Cont. Obs. 1D-Power Plant** along with two belief points represented as uniform

Algorithm 2: $\text{GenRelObs}(\Gamma^{h-1}, a, \mathbf{b}_i) \rightarrow \langle \mathcal{O}^h, p(\mathcal{O}^h | \mathbf{d}\mathbf{x}'_s, a) \rangle$

```

1 begin
2   foreach  $\alpha_j(\mathbf{d}\mathbf{x}'_s) \in \Gamma^{h-1}$  and  $a \in A$  do
3     // Perform exact 1-step DP backup of  $\alpha$ -functions at horizon  $h - 1$ 
4      $\alpha_j^a(\mathbf{d}\mathbf{x}_s, \mathbf{d}\mathbf{x}_o) := \int_{\mathbf{x}'_s} \bigoplus_{\mathbf{d}'_s} p(\mathbf{d}\mathbf{x}_o | \mathbf{d}\mathbf{x}'_s, a) \otimes p(\mathbf{d}\mathbf{x}'_s | \mathbf{d}\mathbf{x}_s, a) \otimes \alpha_j(\mathbf{d}\mathbf{x}'_s) d\mathbf{x}'_s$ 
5     foreach  $\alpha_j^a(\mathbf{d}\mathbf{x}_s, \mathbf{d}\mathbf{x}_o)$  do
6       // Generate value of each  $\alpha$ -vector at belief point  $\mathbf{b}_i(\mathbf{d}\mathbf{x}_s)$  as a function of observations
7        $\delta_j^a(\mathbf{d}\mathbf{x}_o) := \int_{\mathbf{x}_s} \bigoplus_{\mathbf{d}_s} \mathbf{b}_i(\mathbf{d}\mathbf{x}_s) \otimes \alpha_j^a(\mathbf{d}\mathbf{x}_s, \mathbf{d}\mathbf{x}_o) d\mathbf{x}_s$ 
8       // Using casemax, generate observation partitions relevant to each policy – see text for details
9        $\mathcal{O}^h := \text{extract-partition-constraints}[\text{casemax}(\delta_1^{a_1}(\mathbf{d}\mathbf{x}_o), \delta_1^{a_2}(\mathbf{d}\mathbf{x}_o), \dots, \delta_j^{a_r}(\mathbf{d}\mathbf{x}_o))]$ 
10      foreach  $o_k \in \mathcal{O}^h$  do
11        // Let  $\phi_{o_k}$  be the partition constraints for observation  $o_k \in \mathcal{O}^h$ 
12         $p(\mathcal{O}^h = o_k | \mathbf{d}\mathbf{x}'_s, a) := \int_{\mathbf{x}_o} \bigoplus_{\mathbf{d}_o} p(\mathbf{d}\mathbf{x}_o | \mathbf{d}\mathbf{x}'_s, a) \mathbb{I}[\phi_{o_k}] d\mathbf{x}_o$ 
13      return  $\langle \mathcal{O}^h, p(\mathcal{O}^h | \mathbf{d}\mathbf{x}'_s, a) \rangle$ 
14 end

```

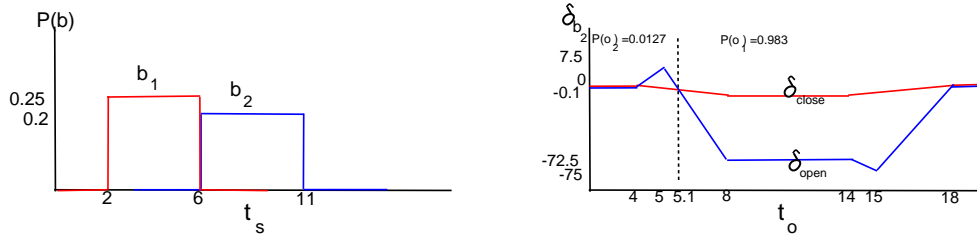


Figure 2: (left) Beliefs b_1, b_2 for **Cont. 1D-Power Plant**; (right) derived observations partitions for b_2 .

distributions: $b_1 : U(t_s; 2, 6)$ and $b_2 : U(t_s; 6, 11)$ as shown in Figure 2 (left). Letting $h = 2$, we assume that Γ^1 contains only the reward as defined in (3). Then in lines 2–4 of Algorithm 2, we compute the following two α -functions for $a \in \{\text{open}, \text{close}\}$ ($j = 1$ since $|\Gamma^1| = 1$):

$$\alpha_1^{\text{close}}(t_s, t_o) = \begin{cases} (t_s < 15) \wedge (t_s - 10 < t_o < t_s) : 10 \\ (t_s \geq 15) \wedge (t_s - 10 < t_o < t_s) : -100 \\ \neg(t_s - 10 < t_o < t_s) : 0 \end{cases} \quad \alpha_1^{\text{open}}(t_s, t_o) = \begin{cases} (t_s - 10 < t_o < t_s) : -0.1 \\ \neg(t_s - 10 < t_o < t_s) : 0 \end{cases}$$

We now need the α -vectors as a function of the observation space for a particular belief state, thus next we marginalize out $\mathbf{d}\mathbf{x}_s$ in lines 5–7. The resulting δ -functions are as follows:

$$\delta_1^{\text{close}}(t_o) = \begin{cases} (14 < t_o < 18) : 0.025t_o - 0.45 \\ (8 < t_o < 14) : -0.1 \\ (4 < t_o < 8) : -0.025t_o - 0.1 \end{cases} \quad \delta_1^{\text{open}}(t_o) = \begin{cases} (15 < t_o < 18) : 25t_o - 450 \\ (14 < t_o < 15) : -2.5t_o - 37.5 \\ (8 < t_o < 14) : -72.5 \\ (5 < t_o < 8) : -25t_o + 127.5 \\ (4 < t_o < 5) : 2.5t_o - 10 \end{cases}$$

Both $\delta_1^{\text{close}}(t_o)$ and $\delta_1^{\text{open}}(t_o)$ are drawn graphically in Figure 2 (right). Here we see that the only relevant observation is whether $t_o < 5.1$ or not – this is the only distinction that determines whether the policy associated with $\delta_1^{\text{close}}(t_o)$ is higher than $\delta_1^{\text{open}}(t_o)$.

The observation dependent δ -functions divide the observation space into regions which can yield the optimal policy according to the belief state b_2 . According to continuous observation space defined in [2], for a 1-dimensional observation space with two or more functions of the observation, we need to find the optimal boundaries or partitions of the space. In their work, numerical solutions are proposed to find the roots of any two observation dependent function. Instead, here we leverage the symbolic power of the max-operator defined in Section 4.1 to find all the boundary regions that define for what partitionings of the observation each δ -function is optimal. For the two δ -functions above the following partitions of the observation space is derived after taking their maximum in lines 9–11:

$$\max \left(\delta_{close}^{b_2}(t_o), \delta_{open}^{b_2}(t_o) \right) = \begin{cases} o_1 : (14 < t_o < 18) & : 0.025t_o - 0.45 \\ o_1 : (8 < t_o < 14) & : -0.1 \\ o_1 : (5.1 < t_o < 8) & : -0.025t_o - 0.1 \\ o_2 : (5 < t_o < 5.1) & : -25t_o + 127.5 \\ o_2 : (4 < t_o < 5) & : 2.5t_o - 10 \end{cases}$$

Note here that each partitioned is labeled according to the original δ -function where o_1 is the label of the partition coming from $\delta_{close}^{b_2}(t_o)$. These partitions define observation regions which we can now use in a similar fashion to a discrete observation set if only the probability of each of the two distinct observation partitions is found. This is demonstrated visually in Figure 2 (right) for b_2 .

Now with the observation partitions derived, all that remains is to calculate probabilities for these relevant observations conditioned on the belief state. Thus we only need to multiply the indicators of each observation partition in this formula to obtain the probability mass lying in each partition:

$$p(o_k | \mathbf{b}_i) := \int_{x'_s} \int_{x_o} \int_{x_s} \bigoplus_{d_o} \bigoplus_{d_s} \bigoplus_{d'_s} p(\mathbf{dx}_o | \mathbf{dx}'_s, a) p(\mathbf{dx}'_s | \mathbf{dx}_s, a) \mathbf{b}_i \mathbb{I}[\phi_{o_k}] d_{x_o} d_{x_s} d_{x'_s}$$

For our example the non-zero probabilities occur in the following partitions as below:

$$p(o_k | b_2) = \begin{cases} (\phi_1 : 14 < t_o < 18) & : 0.2 \\ (\phi_2 : 8 < t_o < 14) & : 0.6 \\ (\phi_3 : 5.1 < t_o < 8) & : 0.183 \\ (\phi_4 : 5 < t_o < 5.1) & : 0.0002 \\ (\phi_5 : 4 < t_o < 5) & : 0.0125 \end{cases}$$

where each constraint is defined as ϕ_i . For any number of actions and observations the total number of observation partitions depends on the number of belief points. In our example we must have two observation probabilities which is defined as the sum of the probabilities for o_1, o_2 according to the constraints:

$$\begin{aligned} o_1 : \phi_1 \vee \phi_2 \vee \phi_3 &\longrightarrow p(o_1 | b_2) = 0.983 \\ o_2 : \phi_4 \vee \phi_5 &\longrightarrow p(o_2 | b_2) = 0.0127 \end{aligned}$$

Hence we can now use the algorithms and methods of a discrete observation setting using the probabilities of the partitioned observation space in PBVI! We note here that our method applies to a restricted class of piecewise functions which is piecewise linear transitions and piecewise constant reward and belief. The main reason is that the integration operation over continuous states and observations only allows constant or linear constraints (upper or lower bounds) over these variables [8]. Thus although in theory we can apply this approach to any piecewise polynomial function, in practice it is limited by the integration bounds. Next we present some results for 2-dimensional continuous observation spaces.

5 Empirical Results

We evaluated our continuous POMDP solution using XADDs on the **1D-Power Plant** example and another variant of this problem with two variables, described below.³

2D-Power Plant: We consider the more complex model of the power plant similar to [1] where the pressure inside the water tank must be controlled to avoid mixing water into the steam or explosion of the tank. We model the pressure variable p as a partially observable variable from the observation readings of the pressure po . The two actions of increase and decrease are defined based on the change in both the temperature and the pressure. For the increase action we define:

$$p(p'_s | \mathbf{p}_s, inc) = \delta \left[p'_s - \begin{cases} (p + 10 > 20) & : 20 \\ \neg(p + 10 > 20) & : p_s + 10 \end{cases} \right] \quad p(t'_s | \mathbf{t}_s, inc) = \delta [t'_s - (t_s + 10)]$$

There is a high reward for staying within the safe temperature and pressure range since it produces power, else depending on how safe it is to have values higher or lower than the safe range, penalty is defined.

³Full problem specifications and Java code to reproduce these experiments are available online in Google Code: <http://code.google.com/p/cpomdp>.

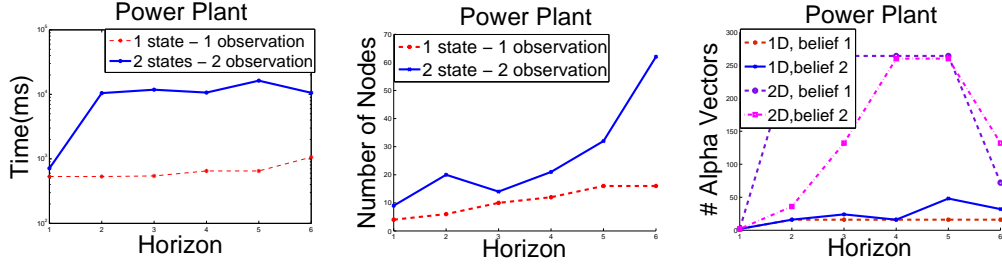


Figure 3: (left) Space vs Horizon; (center) Time vs Horizon; (right) Number of α -vectors vs Horizon.

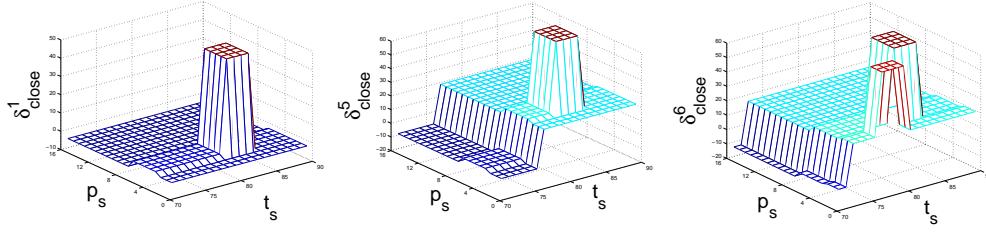


Figure 4: (left) Maximum δ -vector for b_1 and action of *close* in first iteration; (center) $\delta^5_{close}(b_1)$; (right) $\delta^6_{close}(b_1)$. Finer grain partitions show that closing (or opening) the valve can occur at more exact temperatures at higher horizons.

$$R(t_s, p_s, inc) = \begin{cases} (5 \leq p_s \leq 15) \wedge (95 \leq t_s \leq 105) & : 50 \\ (5 \leq p_s \leq 15) \wedge (t_s \leq 95) & : -1 \\ (p_s \geq 15) & : -5 \\ else & : -3 \end{cases}$$

As for the decrease action, the transition functions reduce the temperature by 5 units and the pressure by 10 units as long as the pressure stays above zero. For the reward function, we assume that there is always a small penalty for decreasing the values because power can not be generated. For the observation model we consider two continuous uniform distributions such as the following:

$$p(t_o|t'_s) = \begin{cases} (t_s + 80 < t_o < t_s + 105) & : 0.4 \\ \neg(t_s + 80 < t_o < t_s + 105) & : 0 \end{cases} \quad p(p_o|p'_s) = \begin{cases} (p_s < p_o < p_s + 10) & : 0.1 \\ \neg(p_s < p_o < p_s + 10) & : 0 \end{cases}$$

We define two rectangular uniform beliefs around the regions of rewarding, so that one needs to increase the values while the other should decrease them: $b_1 : U[t_s; 90, 100] * U[p_s; 0, 10]$ and $b_2 : U[t_s; 90, 130] * U[p_s; 10, 30]$. In Figure 3, a time and space analysis of the two versions of **Power Plant** have been performed for up to 6 horizons. As the algorithm progresses, the time required to compute the probability of the partitions and finding the maximum α -vector with respect to beliefs increases for both problem sizes and significantly more for the 2D version. Increase in the problem size, increases the partition numbers on the observation space and this produces more α -vectors which also effects the space required to perform the algorithm. The number of vectors stays the same for most horizons and they drop after convergence in the 2D problem instance. This shows that although the 2D instance takes more time and space than the 1D instance, it still converges within reasonable resources.

In Figure 4 we present plots of the maximum δ -vectors of belief b_1 for different iterations of the 2D problem instance. Starting with the first iteration, the value is highest for the reward range ($5 < p < 15 \wedge 95 < t < 105$) and -1 or less for other places. In the fifth iteration, the value function has partitioned into more pieces, showing how higher temperatures can increase the value without considering the effect of the pressure. In the last plot, horizon $h = 6$ has better tuned the value function so that higher temperatures and pressures increase the value of the maximum δ -vector and also within the reward range, finer grain partitions have been formed.

6 Conclusion

We presented the first exact symbolic operations for PBVI in expressive DC-POMDPs with continuous state *and* observations. Unlike related work that has extended to the continuous state and observation setting [6], we do not approach the problem by sampling. Rather, following [2], the key contribution of this work was to define a discrete set of observation partitions on the multivariate continuous observation space via symbolic maximization techniques and derive the related probabilities using symbolic integration. An important avenue for future work is to determine whether similar techniques can be applied to the difficult case of continuous state, observation, *and* action DC-POMDPs.

Acknowledgments

NICTA is funded by the Australian Government as represented by the Department of Broadband, Communications and the Digital Economy and the Australian Research Council through the ICT Centre of Excellence program. This work was supported by the Fraunhofer ATTRACT fellowship STREAM and by the EC, FP7-248258-First-MM.

References

- [1] Mario Agueda and Pablo Ibarguengoytia. An architecture for planning in uncertain domains. In *Proceedings of the ICTAI 2002 Conference*, Dallas, Texas, 2002.
- [2] Jesse Hoey and Pascal Poupart. Solving pomdps with continuous or large discrete observation spaces. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, Edinburgh, Scotland, 2005.
- [3] Leslie P. Kaelbling, Michael L. Littman, and Anthony R. Cassandra. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101:99–134, 1998.
- [4] G. E. Monahan. Survey of partially observable markov decision processes: Theory, models, and algorithms. *Management Science*, 28(1):1–16, 1982.
- [5] Joelle Pineau, Geoffrey J. Gordon, and Sebastian Thrun. Anytime point-based approximations for large pomdps. *J. Artif. Intell. Res. (JAIR)*, 27:335–380, 2006.
- [6] J. M. Porta, N. Vlassis, M.T.J. Spaan, and P. Poupart. Point-based value iteration for continuous pomdps. *Journal of Machine Learning Research*, 7:195220, 2006.
- [7] Pascal Poupart, Kee-Eung Kim, and Dongho Kim. Closing the gap: Improved bounds on optimal pomdp solutions. In *In Proceedings of the 21st International Conference on Automated Planning and Scheduling (ICAPS-11)*, 2011.
- [8] Scott Sanner and Ehsan Abbasnejad. Symbolic variable elimination for discrete and continuous graphical models. In *In Proceedings of the 26th AAAI Conference on Artificial Intelligence (AAAI-12)*, Toronto, Canada, 2012.
- [9] 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.
- [10] Trey Smith and Reid G. Simmons. Point-based POMDP algorithms: Improved analysis and implementation. In *Proc. Int. Conf. on Uncertainty in Artificial Intelligence (UAI)*, 2005.
- [11] M. Spaan and N. Vlassis. Perseus: Randomized point-based value iteration for pomdps. *Journal of Artificial Intelligence Research (JAIR)*, page 195220, 2005.