



UNIVERSITEIT  
VAN  
AMSTERDAM

IAS technical report IAS-UVA-04-04

## Value Iteration for Continuous-State POMDPs

**Josep M. Porta<sup>†</sup>, Matthijs T.J. Spaan<sup>‡</sup>, and Nikos Vlassis<sup>‡</sup>**

<sup>†</sup>Institut de Robòtica i Informàtica Industrial (UPC-CSIC)  
Llorens i Artigas 4-6, 08028, Barcelona, Spain  
porta@iri.upc.edu

<sup>‡</sup>Informatics Institute, University of Amsterdam  
Kruislaan 403, 1098 SJ Amsterdam, The Netherlands  
{mtjspaan,vlassis}@science.uva.nl

We present a value iteration algorithm for learning to act in Partially Observable Markov Decision Processes (POMDPs) with continuous state spaces. Mainstream POMDP research focuses on the discrete case and this complicates its application to, e.g., robotic problems that are naturally modeled using continuous state spaces. The main difficulty in defining a (belief-based) POMDP in a continuous state space is that expected values over states must be defined using integrals that, in general, cannot be computed in closed form. In this report, we provide three main contributions to the literature on continuous-state POMDPs. First, we show that the optimal finite-horizon value function over the continuous infinite-dimensional POMDP belief space is piecewise linear and convex, and is defined by a finite set of supporting  $\alpha$ -functions that are analogous to the  $\alpha$ -vectors (hyperplanes) defining the value function of a discrete-state POMDP. Second, we show that, for a fairly general class of POMDP models in which all functions of interest are modeled by Gaussian mixtures, all belief updates and value iteration backups can be carried out analytically and exact. Contrary to the discrete case, in a continuous-state POMDP the  $\alpha$ -functions may grow in size (e.g., in the number of Gaussian components) in each value iteration. Third, we show how the recent point-based value iteration algorithms for discrete POMDPs can be extended to the continuous case, allowing for efficient planning in practical problems. In particular, we demonstrate Perseus, our previously proposed randomized point-based value iteration algorithm, in a simple robot planning problem in a continuous domain, where encouraging results are observed.

IAS

intelligent autonomous systems

---

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Preliminaries: POMDPs</b>	<b>1</b>
<b>3</b>	<b>POMDPs in Continuous State Spaces</b>	<b>3</b>
<b>4</b>	<b>Gaussian-based POMDPs</b>	<b>7</b>
4.1	Models for Gaussian-based POMDPs . . . . .	7
4.2	Belief update for Gaussian-Based POMDPs . . . . .	8
4.3	Backup Operator for Gaussian-Based POMDPs . . . . .	9
<b>5</b>	<b>Cs-Perseus: A Point-Based Continuous-State POMDP Solver</b>	<b>10</b>
<b>6</b>	<b>Experiments and Results</b>	<b>13</b>
<b>7</b>	<b>Conclusions and Future Work</b>	<b>15</b>

---

**Intelligent Autonomous Systems**  
Informatics Institute, Faculty of Science  
University of Amsterdam  
Kruislaan 403, 1098 SJ Amsterdam  
The Netherlands  
Tel (fax): +31 20 525 7461 (7490)  
<http://www.science.uva.nl/research/ias/>

**Corresponding author:**  
JM Porta  
tel: +31 20 525 7581  
[porta@science.uva.nl](mailto:porta@science.uva.nl)  
<http://www.science.uva.nl/~porta/>

## 1 Introduction

A popular formalism for decision making under uncertainty is a Markov Decision Process (MDPs) [20]. In this paradigm, an agent interacts with a given system by executing actions, and these actions have the effect of changing the state of the system stochastically, as well as providing rewards/penalties to the agent. The objective of the learning agent is to identify the action that produces the most reward in the long term for each state. When the selection of actions has to be made with uncertain information about the state of the system, the task is naturally formalized as a Partially Observable Markov Decision Process (POMDP) [23, 15, 6, 5, 12]. POMDPs have often been used as a framework for planning in robotics [22, 4, 25, 19]. In general, computing the exact solution of a POMDP is an intractable problem [17, 14], even for the discrete case (i.e., discrete sets of states, actions, and observations). Two main factors cause this high computational cost [18]. The first one is the *curse of history*: the number of action-observation sequences to be considered increases exponentially as we extend the planning horizon. Fortunately the curse of history can be minimized by limiting ourselves to approximate solutions. Recently-developed point-based algorithms [18, 27] are a promising alternative in this line. The second factor that makes POMDP algorithms inefficient is the *curse of dimensionality*: the computational cost of discrete state POMDP algorithms scales with the number of states. Therefore, the finer the granularity of the state space discretization, the higher the cost of solving the POMDP. One insight we can extract from this fact is that it would be desirable to avoid the discretization of the state space. Moreover, real world problems are naturally formalized using continuous spaces. For instance, in a robot navigation problem, the state to be estimated is the pose of the robot that, for a robot moving on a planar surface, is naturally defined in the continuous space of the Cartesian coordinates of the robot and its orientation. Linear POMDPs with continuous states and quadratic reward functions have a closed solution [3]. Existing algorithms for continuous-state POMDPs with general reward functions are based on policy search [16, 1] or approximate (grid-based) value iteration [21, 26]. For discrete-state POMDPs, recent promising algorithms are based on point-based value iteration [18, 27].

In this report, we present a novel approach to solve POMDPs in continuous state spaces via value iteration. The main difficulty of working in continuous state spaces is that expected values over states must be defined using integrals. These integrals cannot be computed in closed form for general functions and, therefore, only approximation techniques can be used [26]. In our approach, we restrict all functions defined on the state space to a particular, although highly expressive, family of functions: linear combinations of Gaussians. This allows us to evaluate all integrals involved in the value iteration POMDP formulation in closed form. Using this fact, we can adapt to the continuous case the rich machinery developed for discrete-state POMDP value iteration, in particular the point-based algorithms.

This report is organized as follows. First, in Section 2, we review the POMDP framework and the value iteration process for discrete-state POMDPs. In Section 3, we generalize the value function representation commonly used in discrete-state POMDPs to continuous-state ones. This allows us to do value iteration for the continuous case. In Section 4, we derive closed formulas for the elements involved in the value iteration framework introduced in Section 3, assuming a Gaussian-based representation for the beliefs and the models defining the POMDP. In Section 5, we use these closed formulas to define a point-based algorithm for Gaussian-based POMDPs. In Section 6, we present some results with the proposed algorithm and, in Section 7, we summarize our work and point to directions for further research.

## 2 Preliminaries: POMDPs

A POMDP models an agent interacting with a system using the following elements

- A set of system states,  $S$ .
- A set of agent actions,  $A$ .
- A set of observations,  $O$ .
- An action (or transition) model defined by  $p(s'|a, s)$ , the probability that the system changes from state  $s$  to  $s'$  when the agent executes action  $a$ .
- An observation model defined by  $p(o|s)$ , the probability that the agent observes  $o$  when the system reaches state  $s$ .
- A reward function defined as  $r_a(s) \in \mathbb{R}$ , the reward obtained by the agent if it executes action  $a$  when the the system is in state  $s$ .

At a given moment, the system is in a state  $s$  and the agent executes an action,  $a$ . As a result, the agent receives a reward  $r$ , the system state changes to  $s'$  and, then, the agent observes  $o$ . The knowledge of the agent about the system state is represented as a *belief*, i.e., a probability distribution over the state space. The initial belief is assumed to be known and, for a discrete set of states, if  $b$  is the belief of the agent about the state, the belief after executing action  $a$  and observing  $o$  is

$$b^{a,o}(s') = \frac{p(o|s')}{p(o|a, b)} \sum_{s \in S} p(s'|s, a) b(s). \quad (1)$$

A function mapping beliefs to actions is called a *policy*. An optimal policy is one that, on the average, generates as much reward as possible. The *value function* condenses the immediate and delayed reward that can be obtained from a given belief. This function is expressed in a recursive way

$$V_n(b) = \max_a Q_n(b, a), \quad (2)$$

with

$$Q_n(b, a) = \sum_{s \in S} r_a(s) b(s) + \gamma \sum_o p(o|b, a) V_{n-1}(b^{a,o}), \quad (3)$$

where  $n$  is the planning horizon,  $S$  and  $O$  are assumed discrete and  $\gamma \in [0, 1)$  is a discount factor that trades off the importance of the immediate and the delayed reward. The above recursion is usually written in functional form

$$V_n = H V_{n-1} \quad (4)$$

and it is known as the *Bellman recursion* [2]. This recursion converges to a fixed point  $V^*$  that is the optimal value function. An optimal policy  $\pi^*$  can be defined as

$$\pi^*(b) = \arg \max_a Q^*(b, a)$$

for  $Q^*$  the  $Q$ -function associated with the optimal value function,  $V^*$ .

*Value iteration* for POMDPs [23, 12, 8] generates a sequence of functions  $V_i$  using the recurrence in Eq. 4 that progressively approach  $V^*$  and computes an approximately optimal policy from the final  $V_i$ .

At first sight the value function seems intractable, but it can be expressed in a simple form [23]

$$V_n(b) = \max_{\{\alpha_n^i\}_i} \sum_s \alpha_n^i(s) b(s),$$

with  $\{\alpha_n^i\}_i$  a set of vectors. Using this formulation, value iteration algorithms typically focus on the computation of the  $\alpha_n$ -vectors.

The initial value function approximation can be the value function with planning horizon 0 [23, 15, 6, 12, 5, 18] or can be defined as a single  $\alpha$ -vector that lower bounds the value function for any possible planning horizon. This second strategy has shown to be more efficient in many benchmark problems

### 3 POMDPs in Continuous State Spaces

In the previous section, we assumed discrete sets of states, actions, and observations. In this section, we generalize POMDPs to continuous state spaces, while still assuming discrete action and observation spaces. With this formulation, we avoid the necessity of discretizing the state space and, thus, we reduce the chance of being affected by the curse of dimensionality.

In the discrete case, expectations for a given belief are computed by summing over the state space (see Eqs. 1 and 3). The generalization to the continuous case amounts to computing these expected values by integrating instead of summing. Thus we have

$$b^{a,o}(s') = \frac{p(o|s')}{p(o|a,b)} \int_s p(s'|s,a) b(s), \quad (5)$$

and

$$Q_n(b,a) = \int_s r_a(s) b(s) + \gamma \sum_o p(o|b,a) V_{n-1}(b^{a,o}), \quad (6)$$

where  $r_a : S \rightarrow \mathbb{R}$  is a continuous reward function for action  $a$ .

With a continuous state space, the belief space is also continuous, as in the discrete case, but now with an infinite number of dimensions. However, there are several properties typical of value functions for discrete state spaces that still hold in the continuous case. Namely, we can prove that (1) the optimal finite-horizon value function is piecewise linear and convex (PWLC) in the belief space, (2) the value function recursion is isotonic, and (3) this recursion is also a contraction (and thus, the iterative computation of the value function for increasing horizons will converge to the optimal value function  $V^*$ ). Next, we prove these three properties

The PWLC is a basic property since it allows to represent the value function using a small set of supporting elements. This kind of representation is the key element to define the value iteration process. To prove this property, we first need to prove the following lemma.

**Lemma 1** *The value function in a continuous-state POMDP can be expressed as*

$$V_n(b) = \max_{\{\alpha_n^i\}_i} \int_s \alpha_n^i(s) b(s), \quad (7)$$

for appropriate  $\alpha$ -functions  $\alpha_n^i : S \rightarrow \mathbb{R}$ .

**Proof:** The proof, as in the discrete case, is done via induction. For planning horizon 0, we only have to take into account the immediate reward and, thus, we have that

$$V_0(b) = \max_a \int_s r_a(s) b(s),$$

and, therefore, if we define

$$\{\alpha_0^i(s)\}_i = \{r_a(s)\}_{a \in A},$$

we have that, as desired

$$V_0(b) = \max_{\{\alpha_0^i\}_i} \int_s \alpha_0^i(s) b(s).$$

For the general case, we have that, using Eqs. 2 and 6,

$$V_n(b) = \max_a \left\{ \int_s r_a(s) b(s) + \gamma \sum_o p(o|a, b) V_{n-1}(b^{a,o}) \right\},$$

and, by the induction hypothesis,

$$V_{n-1}(b^{a,o}) = \max_{\{\alpha_{n-1}^j\}_j} \int_{s'} \alpha_{n-1}^j(s') b^{a,o}(s').$$

From Eq. 5,

$$\begin{aligned} V_{n-1}(b^{a,o}) &= \max_{\{\alpha_{n-1}^j\}_j} \int_{s'} \alpha_{n-1}^j(s') \frac{p(o|s')}{p(o|a, b)} \int_s p(s'|s, a) b(s) \\ &= \frac{1}{p(o|a, b)} \max_{\{\alpha_{n-1}^j\}_j} \int_{s'} \alpha_{n-1}^j(s') p(o|s') \int_s p(s'|s, a) b(s), \end{aligned}$$

and, therefore,

$$\begin{aligned} V_n(b) &= \max_a \left\{ \int_s r_a(s) b(s) + \gamma \sum_o \max_{\{\alpha_{n-1}^j\}_j} \int_{s'} \alpha_{n-1}^j(s') p(o|s') \int_s p(s'|s, a) b(s) \right\} \\ &= \max_a \left\{ \int_s r_a(s) b(s) + \gamma \sum_o \max_{\{\alpha_{n-1}^j\}_j} \int_s \left[ \int_{s'} \alpha_{n-1}^j(s') p(o|s') p(s'|s, a) \right] b(s) \right\}. \end{aligned}$$

At this point, we define

$$\alpha_{a,o}^j(s) = \int_{s'} \alpha_{n-1}^j(s') p(o|s') p(s'|s, a). \quad (8)$$

and we assume that there is a closed form to compute  $\alpha_{a,o}^j$  from  $\alpha_{n-1}^j$  for a given action  $a$  and observation  $o$ . If  $M = |\{\alpha_{n-1}^j\}|$  then, for a given action  $a$  and observation  $o$ , we can generate at most  $M$   $\alpha_{a,o}^j$ -functions. Observe that the  $\alpha_{a,o}^j$  functions are independent of the belief point  $b$  for which we are computing  $V_n$ . With this, we have that

$$V_n(b) = \max_a \left\{ \int_s r_a(s) b(s) + \gamma \sum_o \max_{\{\alpha_{a,o}^j\}_j} \int_s \alpha_{a,o}^j(s) b(s) \right\},$$

and we define

$$\alpha_{a,o,b} = \arg \max_{\{\alpha_{a,o}^j\}_j} \int_s \alpha_{a,o}^j(s) b(s). \quad (9)$$

Observe that, for a given  $a$  and  $o$ ,  $\alpha_{a,o,b}$  is just one of the  $M$  elements in the set  $\{\alpha_{a,o}^j\}_j$ . Using a reasoning parallel to that of the enumeration phase of the Monahan's algorithm [15], we can have, at most,  $|A|M^{|O|}$  different  $\alpha_{a,o,b}$ -functions. The finite cardinality of this set is a crucial point since it proves that we can represent  $V_n(b)$  with a finite set of supporting  $\alpha$ -functions, despite the infinite dimensionality of the belief space.

Using the above, we can write

$$\begin{aligned} V_n(b) &= \max_a \left\{ \int_s r_a(s) b(s) + \gamma \sum_o \int_s \alpha_{a,o,b}(s) b(s) \right\} \\ &= \max_a \left\{ \int_s r_a(s) b(s) + \gamma \int_s \sum_o \alpha_{a,o,b}(s) b(s) \right\} \\ &= \max_a \left\{ \int_s \left[ r_a(s) + \gamma \sum_o \alpha_{a,o,b}(s) \right] b(s) \right\}. \end{aligned}$$

If we define

$$\{\alpha_n^i(s)\}_i = \{r_a(s) + \gamma \sum_o \alpha_{a,o,b}(s)\}_{a \in A}, \quad (10)$$

we have  $V_n$  in the desired form

$$V_n(b) = \max_{\{\alpha_n^i\}_i} \int_s \alpha_n^i(s) b(s), \quad (11)$$

and, thus, the lemma holds.  $\square$

**Lemma 2** *The value function is PWLC in the belief space.*

**Proof:** It holds that

$$V_n(b) = \max_{\{\alpha_n^i\}_i} V_n^i(b),$$

with

$$V_n^i(b) = \int_s \alpha_n^i(s) b(s).$$

For a particular  $V_n^i$  clearly holds

$$V_n^i(\kappa b_1 + \lambda b_2) = \kappa V_n^i(b_1) + \lambda V_n^i(b_2),$$

for arbitrary  $\kappa$  and  $\lambda$ . Therefore, each  $V_n^i$  is a linear function in  $b$ .

The *piecewise linearity* part of the property is given by the fact that the  $\{\alpha_n^i\}_i$  set is of finite cardinality and, as shown above,  $V_n$  is linear, for each individual  $\alpha_n^i$ . Finally, the convexity is given by the fact that we take the maximum of convex (linear) functions when computing the value function and, thus, we obtain a convex function as a result.  $\square$

Eqs. 8 to 10 constitute the value iteration process for continuous state POMDP since they provide a constructive way to determine the elements (i.e., the  $\alpha$ -functions) defining  $V_n$  from those defining  $V_{n-1}$ .

**Lemma 3** *The mapping  $H$  in the value function iteration for continuous state spaces is isotonic.*

**Proof:**  $H$  is said to be isotonic if

$$V \leq U \Rightarrow HV \leq HU.$$

The  $H$  mapping can be seen as

$$HV(b) = \max_a H^a V(b),$$

with

$$H^a V(b) = \int_s r_a(s) b(s) + \gamma \sum_o p(o|a, b) V(b^{a,o}).$$

Let's denote as  $a_1$  the action that maximizes  $HV$  at point  $b$  and  $a_2$  the action that does so for  $HU$

$$\begin{aligned} HV(b) &= H^{a_1} V(b), \\ HU(b) &= H^{a_2} U(b). \end{aligned}$$

By definition, the value for action  $a_1$  for  $HU$  at  $b$  is lower (or equal) than that for  $a_2$ , that is

$$H^{a_1} U(b) \leq H^{a_2} U(b).$$

From a given  $b$  we can compute  $b^{a_1, o}$ , for an arbitrary  $o$  and, then, the following holds

$$\begin{aligned}
V &\leq U \Rightarrow \\
\forall b, o, \quad V(b^{a_1, o}) &\leq U(b^{a_1, o}) \Rightarrow \\
\gamma \sum_o p(o|a_1, b) V(b^{a_1, o}) &\leq \gamma \sum_o p(o|a_1, b) U(b^{a_1, o}) \Rightarrow \\
H^{a_1} V(b) &\leq H^{a_1} U(b) \Rightarrow \\
H^{a_1} V(b) &\leq H^{a_2} U(b) \Rightarrow \\
HV(b) &\leq HU(b) \Rightarrow \\
HV &\leq HU.
\end{aligned}$$

Since  $b$  and, from it  $b^{a_1, o}$ , can be chosen arbitrarily, the value function is isotonic.  $\square$

**Lemma 4** *The mapping  $H$  in the value function iteration for continuous state spaces is a contraction.*

**Proof:** The mapping  $H$  on the value function is said to be a contraction if

$$\|HV - HU\| \leq \beta \|V - U\|,$$

with  $0 \leq \beta < 1$  and  $\|\cdot\|$  the supreme norm. Assume that  $\|HV - HU\|$  is maximum at point  $b$ .  $a_1$  is the optimal action for  $HV$  at  $b$  and so is  $a_2$  for  $HU$

$$\begin{aligned}
HV(b) &= H^{a_1} V(b), \\
HU(b) &= H^{a_2} U(b).
\end{aligned}$$

Then it holds

$$\|HV(b) - HU(b)\| = H^{a_1} V(b) - H^{a_2} U(b).$$

assuming, without loss of generality that  $HV(b) \leq HU(b)$ . Since  $a_1$  is the action that maximizes  $HV$  at  $b$  we have that

$$H^{a_2} V(b) \leq H^{a_1} V(b).$$

Therefore, we have that

$$\begin{aligned}
\|HV - HU\| &= \\
\|HV(b) - HU(b)\| &= \\
H^{a_1} V(b) - H^{a_2} U(b) &\leq \\
H^{a_2} V(b) - H^{a_2} U(b) &= \\
\gamma \sum_o p(o|a_2, b) [V(b^{a_2, o}) - U(b^{a_2, o})] &\leq \\
\gamma \sum_o p(o|a_2, b) \|V - U\| &= \\
\gamma \|V - U\|. &
\end{aligned}$$

Since  $\gamma$  is in  $[0, 1)$ , the lemma holds.  $\square$

Since the value iteration mapping for the continuous state space is a contraction it can be proved that exact value iteration converges in norm to the unique fixed point of  $V = HV$  that is the optimal value function  $V^*$  and that, from the approximations to  $V^*$ , we can derive (near) optimal policies (see [20] Theorems 6.3.1 and 6.2.3).



## 4 Gaussian-based POMDPs

In previous section, we left as an open point how to actually compute the belief update (Eq. 5), the steps in the value iteration process (Eqs. 8 to 10), and the value for a given belief point (Eq. 11). In this section, we show how these computations are possible assuming that the beliefs as well as the observation, action, and reward models are represented as linear combinations of Gaussians. We first formally introduce our assumptions on the models (Section 4.1) and then we define the belief update (Section 4.2) and the basic value iteration steps (Section 4.3) for Gaussian-based POMDPs.

Note that other families of integrable functions could be used to determine the  $\alpha$ -functions in closed form, but Gaussian-based models provide a high degree of flexibility and are of common use in many applications, including robotics [13, 10].

### 4.1 Models for Gaussian-based POMDPs

We will assume that belief points are represented as Gaussian mixtures

$$b(s) = \sum_j w_j \phi(s|s_j, \Sigma_j), \quad (12)$$

with  $\phi$  a Gaussian with mean  $s_j$  and covariance matrix  $\Sigma_j$  and where the mixing weights satisfy  $w_j > 0$ ,  $\sum_j w_j = 1$ . In the extreme case, Gaussian mixtures with an infinite number of components would be necessary to represent a given point in the continuous, infinite-dimensional belief space. However, only Gaussian mixtures with few components are needed in practical situations.

We assume that our observation model is defined non-parametrically from a set of samples  $T = \{(s_i, o_i) \mid i \in [1, N]\}$  with  $o_i$  an observation obtained at state  $s_i$ . Using these samples, the observation model can be defined as

$$p(o|s) = \frac{p(s|o) p(o)}{p(s)},$$

and, assuming a uniform  $p(s)$  in the space covered by  $T$ , and approximating  $p(o)$  from the samples in the training set we have

$$p(o|s) \approx \left[ \frac{1}{N_o} \sum_{i=1}^{N_o} \lambda_i^o \phi(s|s_i^o, \Sigma_i^o) \right] \frac{N_o}{N} = \sum_{i=1}^{N_o} w_i^o \phi(s|s_i^o, \Sigma_i^o)$$

with  $s_i^o$  one of the  $N_o$  points in  $T$  with  $o$  as an associated observation and where  $w_i^o = \lambda_i^o/N$  and  $\Sigma_i^o$  are, respectively, a weighting factor and a covariance matrix associated with that training point. The sets  $\{\lambda_i^o\}_i$  and  $\{\Sigma_i^o\}_i$  should be defined so that

$$p(s) = \sum_o p(s|o) p(o) = \sum_o \sum_{i=1}^{N_o} w_i^o \phi(s|s_i^o, \Sigma_i^o),$$

is (approximately) uniform in the area covered by  $T$  or, in other words, so that

$$\sum_o p(o|s) \approx 1.$$

As far as the action model is concerned, we assume it is linear-Gaussian

$$p(s'|s, a) = \phi(s'|s + \Delta(a), \Sigma^a). \quad (13)$$

with  $\phi$  a Gaussian centered at  $s + \Delta(a)$  with covariance  $\Sigma^a$ . Non-linear action models can be approximated as it is done, for instance, in the extended Kalman filter or in the unscented Kalman filter [11]. The function  $\Delta : A \rightarrow S$  implements the transition model of the system.

Finally, the reward can be seen as an observation with an associated scalar value. Therefore, assuming a finite set of possible rewards  $R = \{r_i \mid i \in [1, M]\}$ , the reward model  $p(r|s, a)$  for each particular  $a$  can be represented in the same way as the observation model

$$p(r|s, a) \approx \sum_{i=1}^{M_r} w_i^r \phi(s|s_i^r, \Sigma_i^r).$$

With that, we have that

$$r_a(s) = \sum_{r \in R} r p(r|s, a) \approx \sum_{r \in R} r \sum_{i=1}^{M_r} w_i^r \phi(s|s_i^r, \Sigma_i^r),$$

that is an unnormalized Gaussian mixture.

## 4.2 Belief update for Gaussian-Based POMDPs

The belief update on Eq. 5 can be implemented in our model taking into account that it consists of two steps. The first one is the application of the action model on the current belief state. This can be computed as the convolution of the Gaussians representing  $b(s)$  (Eq. 12) with the Gaussian representing the action model (Eq. 13). This convolution results in

$$\int_s p(s'|s, a) b(s) = \sum_j w_j \phi(s|s_j + \Delta(a), \Sigma_j + \Sigma^a).$$

In the second step of the belief update, the prediction obtained with the action model is corrected using the information provided by the observation model

$$\begin{aligned} b^{a,o}(s') &\propto \left[ \sum_i w_i^o \phi(s'|s_i^o, \Sigma_i^o) \right] \left[ \sum_j w_j \phi(s|s_j + \Delta(a), \Sigma_j + \Sigma^a) \right] \\ &= \sum_{i,j} w_i^o w_j \phi(s'|s_i^o, \Sigma_i^o) \phi(s|s_j + \Delta(a), \Sigma_j + \Sigma^a) \end{aligned}$$

The product of two Gaussian functions is a scaled Gaussian. Therefore, we have that

$$b^{a,o}(s') \propto \sum_{i,j} w_i^o w_j \delta_{i,j}^{a,o} \phi(s'|s_{i,j}^{a,o}, \Sigma_{i,j}^{a,o}),$$

with

$$\begin{aligned} \delta_{i,j}^{a,o} &= \phi(s_j + \Delta(a) \mid s_i^o, \Sigma_i^o + \Sigma_j + \Sigma^a), \\ \Sigma_{i,j}^{a,o} &= ((\Sigma_i^o)^{-1} + (\Sigma_j + \Sigma^a)^{-1})^{-1}, \\ s_{i,j}^{a,o} &= \Sigma_{i,j}^{a,o} ((\Sigma_i^o)^{-1} s_i^o + (\Sigma_j + \Sigma^a)^{-1} (s_j + \Delta(a))). \end{aligned}$$

Finally, we can re-arrange the terms to get

$$b^{a,o}(s') \propto \sum_k w_k \phi(s'|s_k, \Sigma_k).$$

The proportionality in the definition of  $b^{a,o}(s')$  implies that the weights ( $w_k, \forall k$ ) should be scaled to sum to one.

### 4.3 Backup Operator for Gaussian-Based POMDPs

The computation of the mapping  $H$  (Eq. 4) for a given belief point  $b$  is called a *backup*. This mapping determines the  $\alpha$  function (or  $\alpha$ -vectors in the discrete case) to be included in  $V_n$  for a belief point under consideration (see Eqs. 8 to 10). A full backup, i.e., a backup for the whole belief space, involves the computation of all relevant  $\alpha$ -functions for  $V_n$ . Full backups are computationally expensive (in the discrete case they involve the use of linear programming in order to determine a sufficient set of points on which to backup), but the backup for a single belief point is relatively cheap. This is exploited by the point-based POMDP algorithms to efficiently approximate  $V_n$  on a fixed set of belief points [18, 27]. Next, we describe the backup operator on a continuous state space that we will use later in the PERSEUS algorithm.

The backup for a given belief point  $b$  is

$$\text{backup}(b) = \arg \max_{\{\alpha_n^i\}_i} \int_s \alpha_n^i(s) b(s),$$

where  $\alpha_n^i(s)$  is defined in Eqs. 9 and 10 from the  $\alpha_{a,o}$ -functions (Eq. 8).

**Lemma 5** *The functions  $\alpha_n^i(s)$  can be expressed as linear combinations of Gaussians, assuming the sensor, action and reward models are also Gaussian-based.*

**Proof:** This lemma can be proved via induction. For  $n = 0$ ,  $\alpha_0^i(s) = r_a(s)$  for a fixed  $a$  and thus it is indeed an unnormalized Gaussian mixture. For  $n > 0$ , we assume that

$$\alpha_{n-1}^j(s') = \sum_k w_k^j \phi(s' | s_k^j, \Sigma_k^j).$$

Then, with our particular models,  $\alpha_{a,o}^j(s)$  in Eq. 8 is the integral of three linear combinations of Gaussians

$$\begin{aligned} \alpha_{a,o}^j(s) &= \int_{s'} \left[ \sum_k w_k^j \phi(s' | s_k^j, \Sigma_k^j) \right] \left[ \sum_l w_l^o \phi(s' | s_l^o, \Sigma_l^o) \right] \phi(s' | s + \Delta(a), \Sigma^a) \\ &= \int_{s'} \sum_{k,l} w_k^j w_l^o \phi(s' | s_k^j, \Sigma_k^j) \phi(s' | s_l^o, \Sigma_l^o) \phi(s' | s + \Delta(a), \Sigma^a) \\ &= \sum_{k,l} w_k^j w_l^o \int_{s'} \phi(s' | s_k^j, \Sigma_k^j) \phi(s' | s_l^o, \Sigma_l^o) \phi(s' | s + \Delta(a), \Sigma^a). \end{aligned}$$

In this case, we have to perform the product of two Gaussians twice, once for  $\phi(s' | s_k^j, \Sigma_k^j)$  and  $\phi(s' | s_l^o, \Sigma_l^o)$  to get  $(\delta_{k,l}^{j,o} \phi(s' | s_1, \Sigma_1))$  and once more for  $(\delta_{k,l}^{j,o} \phi(s' | s_1, \Sigma_1))$  and  $\phi(s' | s + \Delta(a), \Sigma^a)$  to get  $(\delta_{k,l}^{j,o} \beta_{k,l}^{j,o,a}(s) \phi(s' | s, \Sigma))$ . The terms  $\delta_{k,l}^{j,o}$  and  $\beta_{k,l}^{j,o,a}(s)$  can be expressed as

$$\begin{aligned} \delta_{k,l}^{j,o} &= \phi(s_l^o | s_k^j, \Sigma_k^j + \Sigma_l^o), \\ \beta_{k,l}^{j,o,a}(s) &= \phi(s | s_{k,l}^{j,o} - \Delta(a), \Sigma_{k,l}^{j,o} + \Sigma^a), \end{aligned}$$

with

$$\begin{aligned} \Sigma_{k,l}^{j,o} &= [(\Sigma_k^j)^{-1} + (\Sigma_l^o)^{-1}]^{-1}, \\ s_{k,l}^{j,o} &= \Sigma_{k,l}^{j,o} [(\Sigma_k^j)^{-1} s_k^j + (\Sigma_l^o)^{-1} s_l^o]. \end{aligned}$$

With this, we have

$$\begin{aligned}\alpha_{a,o}^j(s) &= \sum_{k,l} w_k^j w_l^o \int_{s'} \delta_{k,l}^{j,o} \beta_{k,l}^{j,o,a}(s) \phi(s'|s, \Sigma) \\ &= \sum_{k,l} w_k^j w_l^o \delta_{k,l}^{j,o} \beta_{k,l}^{j,o,a}(s) \int_{s'} \phi(s'|s, \Sigma) \\ &= \sum_{k,l} w_k^j w_l^o \delta_{k,l}^{j,o} \beta_{k,l}^{j,o,a}(s).\end{aligned}$$

Once we have the  $\alpha_{a,o}^j$ -functions, we can compute the  $\alpha_n^i$ -functions. To do that, we need to determine the  $\alpha_{a,o}^j$  for which

$$\int_s \alpha_{a,o}^j(s) b(s)$$

is maximized. Since the integral of the product of two Gaussian mixtures (in particular an  $\alpha$ -function and a belief point) is a rather common operation in the continuous state POMDP framework we will denote it by

$$\langle \alpha, b \rangle = \int_s \alpha(s) b(s).$$

This operator can be computed as

$$\begin{aligned}\langle \alpha, b \rangle &= \int_s \left[ \sum_k w_k \phi(s|s_k, \Sigma_k) \right] \left[ \sum_l w_l \phi(s|s_l, \Sigma_l) \right] \\ &= \sum_{k,l} w_k w_l \int_s \phi(s|s_k, \Sigma_k) \phi(s|s_l, \Sigma_l) \\ &= \sum_{k,l} w_k w_l \phi(s_l|s_k, \Sigma_k + \Sigma_l) \int_s \phi(s|s_{k,l}, \Sigma_{k,l}) \\ &= \sum_{k,l} w_k w_l \phi(s_l|s_k, \Sigma_k + \Sigma_l).\end{aligned}$$

Using this operator and Eqs. 9 and 10, we define

$$\{\alpha_n^i(s)\}_i = \{r_a(s) + \gamma \sum_o \arg \max_{\{\alpha_{a,o}^j\}_j} \langle \alpha_{a,o}^j, b \rangle\}_{a \in A}.$$

Since all elements involved in the definition are linear combination of Gaussians so is the final result.  $\square$  Using the above lemma, the **backup** function is

$$\text{backup}(b) = \arg \max_{\{\alpha_n^i\}_i} \langle \alpha_n^i, b \rangle,$$

and the value of  $V_n$  at  $b$  (Eq. 11) is simply

$$V_n(b) = \langle \text{backup}(b), b \rangle.$$

## 5 Cs-Perseus: A Point-Based Continuous-State POMDP Solver

In this section, we use the **backup** operator to define a point-based approximate continuous-state POMDP solver. In particular, we show how to extend to the continuous case the point-based value iteration algorithm PERSEUS [27, 24], which has been shown to be very efficient for discrete

```

Perseus
Input: A continuous state POMDP.
Output:  $V_n$ , an approximation to the optimal
          value function,  $V^*$ .

1: Initialize
2:    $B \leftarrow$  A set of randomly sampled belief points.
3:    $\alpha \leftarrow \frac{\min\{R\}}{1-\gamma} \phi(s|0, \Sigma_\infty)$ 
4:    $n \leftarrow 0$ 
5:    $V_n \leftarrow \{\alpha\}$ 
6:   do
7:      $\forall b \in B,$ 
8:        $\text{Function}_n(b) \leftarrow \arg \max_{\alpha \in V_n} \langle \alpha, b \rangle$ 
9:        $\text{Value}_n(b) \leftarrow \langle \text{Function}_n(b), b \rangle$ 
10:     $V_{n+1} \leftarrow \emptyset$ 
11:     $\tilde{B} \leftarrow B$ 
12:    do
13:       $b \leftarrow$  Point sampled randomly from  $\tilde{B}$ .
14:       $\alpha \leftarrow \text{backup}(b)$ 
15:      if  $\langle \alpha, b \rangle < \text{Value}_n(b)$ 
16:         $\alpha \leftarrow \text{Function}_n(b)$ 
17:         $\tilde{B} \leftarrow \tilde{B} \setminus \{b' \in \tilde{B} \mid \text{Function}_n(b') = \alpha\}$ 
18:      else
19:         $\tilde{B} \leftarrow \tilde{B} \setminus \{b' \in \tilde{B} \mid \langle \alpha, b' \rangle \geq \text{Value}_n(b')\}$ 
20:      endif
21:       $V_{n+1} \leftarrow V_{n+1} \cup \{\alpha\}$ 
22:    until  $\tilde{B} = \emptyset$ 
23:     $n \leftarrow n + 1$ 
24: until convergence

```

**Table 1:** The PERSEUS algorithm. The **backup** function is described in Section 4.3.

state POMDPs. The continuous-state PERSEUS algorithm is shown in Table 1. Point-based POMDP algorithms focus on identifying the  $\alpha$ -functions ( $\alpha$ -vectors in the discrete case) for the belief points where the agent is more likely to be. The  $\alpha$ -functions for this restricted set of belief points generalize over the whole belief space and, thus, they can be used to approximate the value function for any belief point. The result is an approximation of the value function with less error in regions of the belief space where decisions are more likely to be taken.

The value update scheme of PERSEUS implements a randomized approximate value function recursion  $V_n = \tilde{H}V_{n-1}$  for a set of randomly sampled belief points  $B$ . First (Table 1, line 2), we let the agent randomly explore the environment and collect a set  $B$  of reachable belief points. Next (Table 1, lines 3-5), we initialize the value function  $V_0$  as a single weighted Gaussian with large covariance and with weight  $\min\{R\}/(1-\gamma)$ , with  $R$  the set of possible rewards.

Starting with  $V_0$ , PERSEUS performs a number of approximate value function update stages. The definition of the value update process can be seen on lines 10–22 in Table 1, where  $\tilde{B}$  is a set of non-improved points: points for which  $V_{n+1}(b)$  is still lower than  $V_n(b)$ . At the start of each update stage,  $V_{n+1}$  is set to  $\emptyset$  and  $\tilde{B}$  is initialized to  $B$ . As long as  $\tilde{B}$  is not empty, we sample a point  $b$  from  $\tilde{B}$  and compute the new  $\alpha$ -function associated with this point using the **backup** operator (see Section 4.3). If this  $\alpha$ -function improves the value of  $b$  (i.e., if  $\langle \alpha, b \rangle \geq V_n(b)$ ), we add  $\alpha$  to  $V_{n+1}$ . The hope is that  $\alpha$  improves the value of many other points, and all these points

**Gaussian Mixture Condensation(f, m)****Input:** A Gaussian mixture  $f = \sum_{i=1}^k w_i f_i(x|\mu_i, \Sigma_i)$ .The maximum number of components  
in the output mixture,  $m, m < k$ .**Output:** A Gaussian mixture  $g = \sum_{i=1}^m w'_i g_i(x|\mu'_i, \Sigma'_i)$  that  
locally minimizes  $\sum_{i=1}^k w_i \min_{j \in [1, m]} KL(f_i \| g_j)$ 1: **Initialize**2:     **for**  $j = 1$  **to**  $m$ 3:          $w'_j \leftarrow w_j$ 4:          $\mu'_j \leftarrow \mu_j$ 5:          $\Sigma'_j \leftarrow \Sigma_j$ 6:      $d \leftarrow \sum_{i=1}^k w_i \min_{j \in [1, m]} KL(f_i \| g_j)$ 7:     **do**8:         **Compute the mapping from**  $f$  **to**  $g$ 9:         **for**  $i = 1$  **to**  $k$ 10:              $\pi(i) \leftarrow \arg \min_{j \in [1, m], w'_j > 0} KL(f_i \| g_j)$ 11:         **Define a new**  $g$ 12:         **for**  $j = 1$  **to**  $m$ 13:              $I_j \leftarrow \{i \mid \pi(i) = j, i \in [1, k]\}$ 14:              $w'_j \leftarrow \sum_{i \in I_j} w_i$ 15:              $\mu'_j \leftarrow \frac{1}{w'_j} \sum_{i \in I_j} w_i \mu_i$ 16:              $\Sigma'_j \leftarrow \frac{1}{w'_j} \sum_{i \in I_j} w_i (\Sigma_i + (\mu_i - \mu'_j)(\mu_i - \mu'_j)^\top)$ 17:      $d' \leftarrow d$ 18:      $d \leftarrow \sum_{i=1}^k w_i KL(f_i \| g_{\pi(i)})$ 19: **until**  $\frac{|d-d'|}{d} < \epsilon$ **Table 2:** Gaussian mixture condensation algorithm.  $\epsilon$  is a sufficiently small threshold.

are removed from  $\tilde{B}$ . Often, a small number of vectors will be sufficient to improve  $V_n(b) \forall b \in B$ , especially in the first steps of value iteration. As long as  $\tilde{B}$  is not empty we continue sampling belief points from it and trying to add their  $\alpha$ -functions to  $V_{n+1}$ .

If the  $\alpha$  computed by the **backup** operator does not improve at least the value of  $b$  (i.e.,  $\langle \alpha, b \rangle < V_n(b)$ , see lines 15–17 in Table 1), we ignore  $\alpha$  and insert a copy of the maximizing function of  $b$  from  $V_n$  in  $V_{n+1}$ . Point  $b$  is now considered improved and is removed from  $\tilde{B}$ , together with any other belief points that had the same function as maximizing one in  $V_n$ . This procedure ensures that  $\tilde{B}$  shrinks at each iteration and that the value update stage terminates.

PERSEUS stops when a given convergence criterion holds. This criterion can be based on the stability of the value function, on the stability of the associated policy, or simply on a maximum number of iterations.

One point that deserves special consideration when implementing the PERSEUS algorithm is the possible explosion of the number of components in the Gaussian mixtures defining the  $\alpha$ -functions for increasing  $n$ 's and on the number of components in the belief representation when the belief update (see Section 4.2) is repeated for many time steps. If  $C_o$  is the average number of components in the observation model and  $C_b$  is the average number of components in the belief, the number of components in the  $\alpha$ -functions or in the belief after  $n$  iterations scales with  $O(C_o^n C_b^n)$ . Since the larger the number of components the slower the basic operations of the algorithm and efficient implementation of the algorithm requires to keep the number of

components reasonably bounded. To achieve this objective, we use the procedure described in [7] that transforms a given Gaussian mixture with  $k$  components to another Gaussian mixture with at most  $m$  components,  $m < k$ , while retaining the initial component structure. The algorithm is detailed in Table 2.

The algorithm uses the Kullback-Leibler,  $KL$ , distance between two Gaussian distributions  $f_i = N(\mu, \Sigma)$ ,  $g_j = N(\mu', \Sigma')$  that is

$$KL(f_i \| g_j) = \frac{1}{2} \left( \log \frac{|\Sigma'|}{|\Sigma|} + \text{Tr}((\Sigma')^{-1}\Sigma) + (\mu - \mu')^\top (\Sigma')^{-1} (\mu - \mu') - c \right)$$

with  $c$  the dimensionality of the space where the Gaussians are defined.

Observe that the above procedure is defined for normalized Gaussian mixtures and our  $\alpha$  functions are unnormalized Gaussian mixtures. Therefore, for the  $\alpha$ -function compression, we use a modified version of the procedure just described where the weights are normalized after taking its absolute value (so that relevant reward peaks either negative or positive are preserved). After the compression, the reverse procedure is used to recover weights in the original scale.

In our implementation, we limit the number of components in the  $\alpha$  functions to those in the  $\alpha$  functions in  $V_0$ .

A similar number of components explosion occurs when computing the belief update detailed in Section 4.2. In this case, we use the Gaussian mixtures clustering algorithm so that number of components never exceeds that of the initial belief.

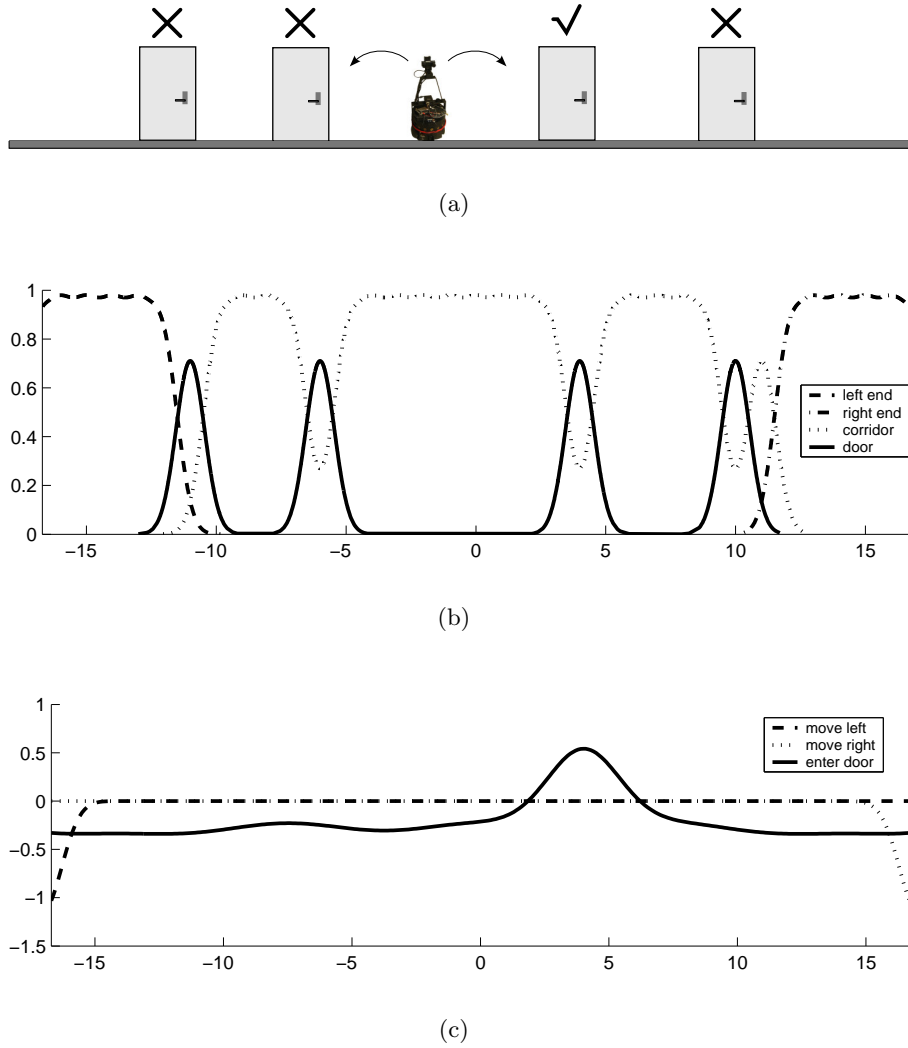
## 6 Experiments and Results

To demonstrate the viability of our method we carried out an experiment in a robotic domain. In this problem (see Fig. 1-a), a robot is moving in a corridor with four doors. The robot can detect when it is in front of a door and when it is at the left or right end of the corridor. In any other situation, the robot just detects that it is in a corridor (see Fig. 1-b). The robot can move 2 units to the left or to the right (with  $\Sigma^a = 0.05$ ) and can try to enter a door at any point (even when not in front of a door). The target for the robot is to locate the second door from the right and to enter it. The robot only gets positive reward when it enters the target door (see Fig. 1-c). When the robot tries to move further than the end of the corridor (either at the right or at the left) or when it tries to enter the door at a wrong position it gets negative reward.

The set of beliefs  $B$  used in the CS-PERSEUS algorithm contains 1000 unique belief points. Those belief points are collected using random walks departing from a belief including 4 components that approximate a uniform distribution on the whole corridor. The walks of the robot along the corridor are organized in episodes where the robot executes actions until it tries to enter a door or until it executes 25 (movement) actions.

The experimental setup is completed by setting  $\gamma$  to 0.95, compressing beliefs so that they never contain more than 4 components (i.e., the number of components of the initial belief) and compressing  $\alpha$ -functions so that they never have more components than those used to represent the reward function (11 components).

Fig. 2 shows the average results obtained after 10 runs of the CS-PERSEUS algorithm on this problem. The first plot (top-left) shows that the value computed as  $\sum_b V(b)$  converges. The second plot (top-right) shows the expected discounted reward averaged for 100 episodes with the policy available at the corresponding time slice. The plot indicates that the robot successfully learns to find out its position and to distinguish between the four doors. Next plot (bottom-left) shows the number of  $\alpha$ -functions used to represent the value function. We can see that the number of  $\alpha$ -functions used increases, but is far below 1000, the maximum possible



**Figure 1:** A pictorial representation of the test problem (a), the corresponding observation model (b) and the reward model (c).

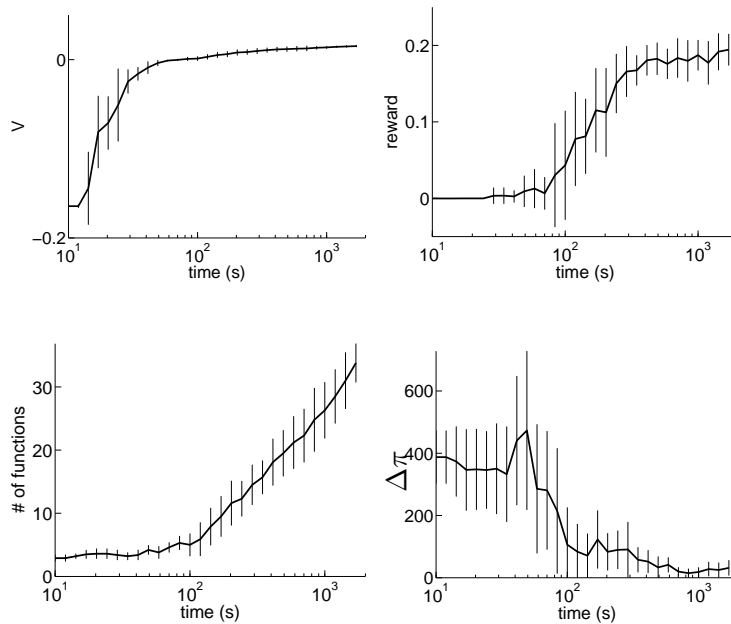
number of  $\alpha$ -functions (in the extreme case we would use a different  $\alpha$ -function for each point in  $B$ ). Finally (plot at Fig. 2, bottom-right) we show the number of changes in the policy from one time step to the next one. The changes in the policy are computed as the number of elements in  $B$  with a different action from one time slice to the next. The number of policy changes drop to close to zero, indicating convergence with respect to the particular  $B$ .

Following the learned policy the robot moves to one of the ends of the corridor to determine its position and then towards the correct door to enter it. Fig. 3 shows the evolution of the belief of the robot and the executed action in each case from the initial stage of the episode to the point at which the target door is reached.

Finally, Fig. 4 plots the value for beliefs with only one component parametrized by the average and the covariance of this component. We can see that, as the uncertainty about the position of the robot grows (i.e., as the covariance is larger) the value of the corresponding belief decreases. The colors in the figure correspond to the different actions: light-gray for moving to the right, white for entering the door, and dark-gray for moving to the left.

Observe that the advantage of using a continuous state space is that we obtain a scale-invariant solution. If we have to solve the same problem in a longer corridor, we can just scale





**Figure 2:** Top: Evolution of the value for all the beliefs in  $B$  and the average accumulated discounted reward for 100 episodes. Bottom: Number of vectors in  $V_n$  and the number of policy changes. Results are averaged for 10 repetitions and the bars represent the standard deviation.

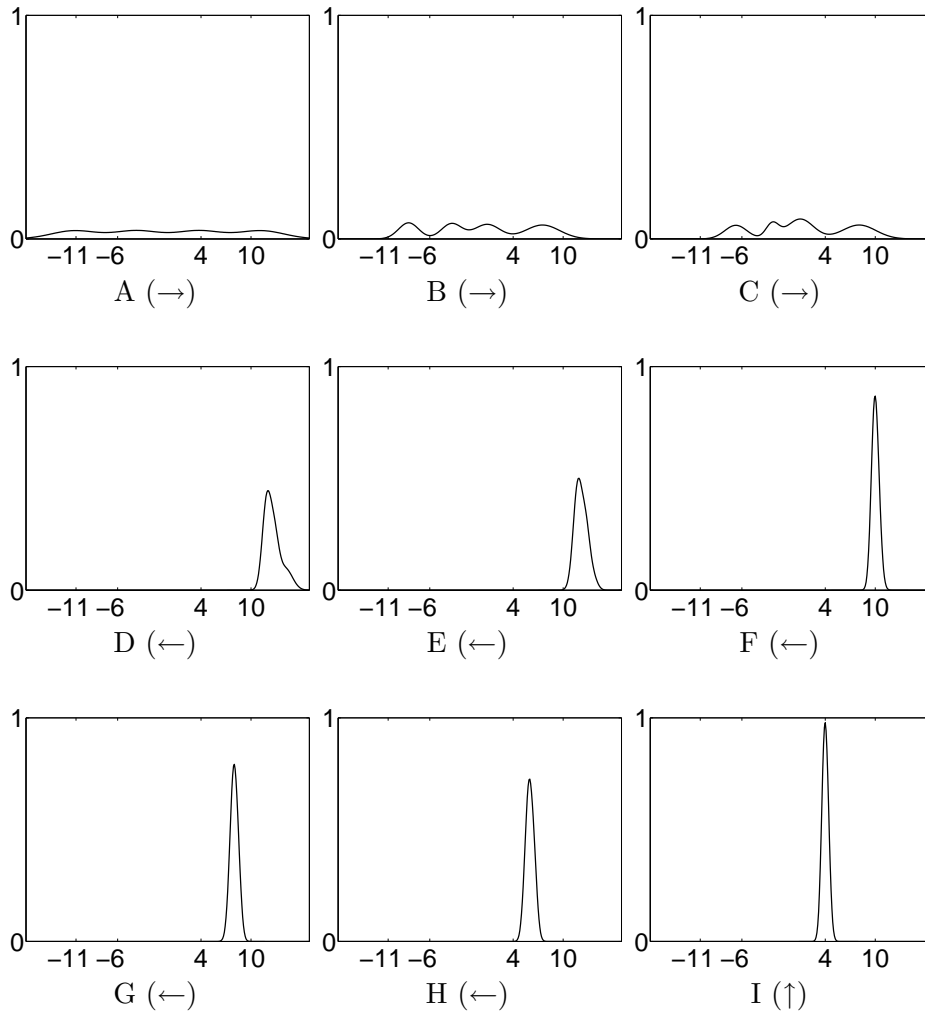
the Gaussians used in the problem definition and we will obtain the solution with the same cost as we have now. The only difference is that more actions would be needed in each episode to reach the correct door. When discretizing the environment, the granularity has to be in accordance with the size of the actions taken by the robot ( $\pm 2$  left/right) and, thus, the number of states and, consequently, the cost of the planning grow as the environment grows.

## 7 Conclusions and Future Work

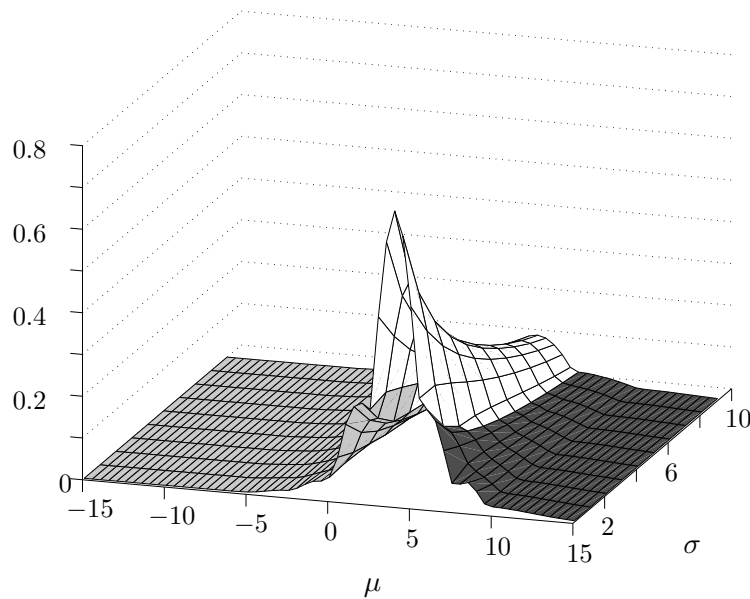
In this paper we have shown how to generalize value iteration to continuous-state POMDPs, and in particular for the case of Gaussian-based beliefs and models. This allowed us to define an efficient point-based value iteration algorithm that seems to be appropriate for planning problems that are often encountered in robotics.

An approach to continuous-state POMDPs that is closely related to ours is presented in [26]. In that work, a belief is represented by a set of weighted samples, which can be regarded as a degenerate version of our Gaussian mixture representation. Additionally, the value function is approximated by nearest-neighbor interpolation, whereas in our case the value function achieves generalization through a set of  $\alpha$ -functions. Also, in the above work a real-time dynamic programming approach is used for updating the value function, with the Bellman backup operator being approximated by sampling from the belief transition model. In our case, value iteration applies on a pre-collected set of beliefs, while the Bellman backup operator is analytically computed given the particular value function representation. Although we have not directly compared our method to the method presented in [26], we expect our method to be faster (since it plans on a fixed set of belief points) and the value function to generalize better over the belief space (through the use of  $\alpha$ -functions).

Ongoing work involves extending our framework to continuous action [24] and observation spaces [9], as well as defining approximate belief representations using Monte Carlo tech-



**Figure 3:** Evolution of the belief when following the discovered policy. The arrows under the snapshots represent the actions:  $\rightarrow$  for moving right,  $\leftarrow$  for moving left and  $\uparrow$  for entering the door. On the  $x$ -axis the four door locations are indicated.



**Figure 4:** Value function for single component beliefs as a function of the average and the covariance.

niques [26].

## References

- [1] D. Aberdeen and J. Baxter. Scalable Internal-State Policy-Gradient Methods for POMDPs. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 3–10, 2002.
- [2] R.E. Bellman. *Dynamic Programming*. Princenton University Press, 1957.
- [3] D.P. Bertsekas. *Dynamic Programming and Optimal Control*. Athena Scientific cop, Belmont, MA, 2 edition, 2001.
- [4] A.R. Cassandra, L.P. Kaelbling, and J.A. Kurien. Acting under Uncertainty: Discrete Bayesian Models for Mobile-Robot Navigation. In *Proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 1996.
- [5] A.R. Cassandra, M.L. Littman, and N.L. Zhang. Incremental Pruning: A Simple, Fast, Exact Algorithm for Partially Observable Markov Decision Processes. In *Proceedings of the Thirteenth Annual Conference on Uncertainty in Artificial Intelligence (UAI-97)*, 1997.
- [6] H.T. Cheng. *Algorithms for Partially Observable Markov Decision Processes*. PhD thesis, University of British Columbia, 1988.
- [7] J. Goldberger and S. Roweis. Hierarchical Clustering of a Mixture Model. In *Advances in Neural Information Processing Systems (NIPS)*, 2005.
- [8] M. Hauskrecht. Value Function Approximations for Partially Observable Markov Decision Processes. *Journal of Artificial Intelligence Research*, 13:33–95, 2000.
- [9] J. Hoey and P. Poupart. Solving POMDPs with Continuous or Large Discrete Observation Spaces. In *Proceedings of the International Joint Conference on Artificial Intelligence*, 2005.

- [10] P. Jensfelt and S. Kristensen. Active Global Localization for a Mobile Robot Using Multiple Hypothesis Tracking. *IEEE Transactions on Robotics and Automation*, 17(5):748–760, 2001.
- [11] J. Julier and J. K. Uhlmann. A New Extension of the Kalman Filter to Nonlinear Systems. In *In Proceedings of AeroSense: The 11th International Symposium on Aerospace/Defence Sensing, Simulation and Controls*, 1997.
- [12] L.P. Kaelbling, M.L. Littman, and A.R. Cassandra. Planning and Acting in Partially Observable Stochastic Domains. *Artificial Intelligence*, 101(1-2):99–134, 1998.
- [13] J.J. Leonard and H.F. Durrant-Whyte. Mobile Robot Localization by Tracking Geometric Beacons. *IEEE Transactions on Robotics and Automation*, 7(3):376–382, 1991.
- [14] O. Madani, S. Hanks, and A. Condon. On the Undecidability of Probabilistic Planning and Infinite-Horizon Partially Observable Markov Decision Problems. In *Proceedings of the Sixteenth National Conference on Artificial Intelligence (AAAI)*, pages 541–548, 1999.
- [15] G.E. Monahan. A Survey of Partially Observable Markov Decision Processes: Theory, Models, and Algorithms. *Management Science*, 28(1):1–16, 1982.
- [16] A.Y. Ng and M. Jordan. PEGASUS: A Policy Search Method for Large MDPs and POMDPs. In *Proceedings of the 16th Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 406–415, 2000.
- [17] C. Papadimitriou and J.N. Tsitsiklis. The Complexity of Markov Decision Processes. *Mathematical and Operations Research*, 12(3):441–450, 1987.
- [18] J. Pineau, G. Gordon, and S. Thrun. Point-based Value Iteration: An Anytime Algorithm for POMDPs. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2003.
- [19] J. Pineau, M. Montemerlo, M. Pollack, N. Roy, and S. Thrun. Towards Robotic Assistants in Nursing Homes: Challenges and Results. In *Robotics and Autonomous Systems*, volume 42, pages 271–281, 2003.
- [20] M.L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley Series in Probability and Mathematical Statistics. John Wiley and Sons, Inc., 1994.
- [21] N. Roy, G. Gordon, and S. Thrun. Finding Approximate POMDP Solutions Through Belief Compression. *Journal of Artificial Intelligence Research*, 23:1–40, 2005.
- [22] R. Simmons and S. Koenig. Probabilistic Robot Navigation in Partially Observable Environments. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, 1995.
- [23] E.J. Sondik. *The Optimal Control of Partially Observable Markov Processes*. PhD thesis, Stanford University, 1971.
- [24] M.T.J. Spaan and N. Vlassis. Perseus: Randomized Point-based Value Iteration for POMDPs. *Journal of Artificial Intelligence Research*, 2005.
- [25] G. Theodorou and S. Mahadevan. Approximate Planning with Hierarchical Partially Observable Markov Decision Processes for Robot Navigation. In *IEEE International Conference on Robotics and Automation*, pages 1347–1352, 2002.
- [26] S. Thrun. Monte Carlo POMDPs. In S.A. Solla, T.K. Leen, and K.-R. Müller, editors, *Advances in Neural Information Processing Systems (NIPS)*, pages 1064–1070. MIT Press, 2000.

- 
- [27] N. Vlassis and M.T.J. Spaan. A Fast Point-Based Algorithm for POMDPs. In *In Proceedings of Annual Machine Learning Conference of Belgium and the Netherlands, Brussels, Belgium*, pages 170–176, 2004.



---

## Acknowledgements

We would like to thank J.J. Verbeek and W. Zajdel for their contributions to the work reported here. This work was developed while J.M. Porta was visiting the IAS Group of the University of Amsterdam.

He has been partially supported by Ramón y Cajal contract from the Spanish Ministry for Science and Technology. M.T.J. Spaan and N. Vlassis are supported by PROGRESS, the embedded systems research program of the Dutch organization for Scientific Research NWO, the Dutch Ministry of Economic Affairs and the Technology Foundation STW, project AES 5414. Authors are listed in alphabetical order.

## IAS reports

This report is in the series of IAS technical reports. The series editor is Stephan ten Hagen ([stephanh@science.uva.nl](mailto:stephanh@science.uva.nl)). Within this series the following titles appeared:

W. Zajdel, A.T. Cemgil and B.J.A. Kröse. *A Hybrid Graphical Model for Online Multi-camera Tracking*. Technical Report IAS-UVA-04-03, Informatics Institute, University of Amsterdam, The Netherlands, November 2004.

Matthijs T.J. Spaan and Nikos Vlassis. *Perseus: Randomized point-based value iteration for POMDPs*. Technical Report IAS-UVA-04-02, Informatics Institute, University of Amsterdam, The Netherlands, November 2004.

Nikos Vlassis and Jakob J Verbeek. *Gaussian mixture learning from noisy data*. Technical Report IAS-UVA-04-01, Informatics Institute, University of Amsterdam, The Netherlands, September 2004.

Jelle R. Kok and Nikos Vlassis. *The Pursuit Domain Package*. Technical Report IAS-UVA-03-03, Informatics Institute, University of Amsterdam, The Netherlands, August 2003.

Joris Portegies Zwart, Ben Kröse, and Sjoerd Gelsema. *Aircraft Classification from Estimated Models of Radar Scattering*. Technical Report IAS-UVA-03-02, Informatics Institute, University of Amsterdam, The Netherlands, January 2003.

Joris Portegies Zwart, René van der Heiden, Sjoerd Gelsema, and Frans Groen. *Fast Translation Invariant Classification of HRR Range Profiles in a Zero Phase Representation*. Technical Report IAS-UVA-03-01, Informatics Institute, University of Amsterdam, The Netherlands, January 2003.

All IAS technical reports are available for download at the IAS website, <http://www.science.uva.nl/research/ias/publications/reports/>.