

Accelerating Point-Based Value Iteration via Active Sampling of Belief Points and Gaussian Process Regression

Siqiong Zhou¹, Ashif S. Iquebal¹, and Esma S. Gel²

¹School of Computing and Augmented Intelligence, Arizona State University, USA

²College of Business, University of Nebraska-Lincoln, USA

Abstract

Partially Observable Markov Decision Processes (POMDPs) are fundamental to decision-making under uncertainty. We introduce a novel scalable approach to accelerate upper bound estimation in Point-Based Value Iteration (PBVI) algorithms, the leading method to solve large-scale POMDPs. PBVI approximates the value function using a set of belief points rather than the entire continuous belief space and relies on lower and upper bounds for convergence. While lower bounds are straightforward to compute, PVBI requires repeated sawtooth projection operations to approximate the upper bound convex hull, significantly increasing the computational burden although many of these sawtooth projections become redundant as the belief set expands. To address this, we infer the upper bound using the upper confidence bound of a Gaussian Process Regression (GP-UCB) fitted over a subset of the most informative reachable belief points—the ones that exhibit linear independence in some high-dimensional Hilbert space. This approach reduces the number of sawtooth projections by 84.3% on average without compromising the solution quality. We further establish the theoretical consistency of the proposed GP-UCB estimate of the upper bound and show convergence to the true upper bound convex hull. We implement GP-UCB and test its performance using five benchmark finite-horizon POMDPs, demonstrating its effectiveness in estimating upper bounds and improving PBVI performance. GP-UCB reduces computation time by 30% to 60% on smaller problems and up to 99.7% on larger ones, while achieving the same gaps as the pure sawtooth projection method.

Keywords: partially observable Markov decision processes, Gaussian process regression, upper confidence bound, point-based value iteration

1 Introduction

A Partially Observable Markov Decision Process (POMDP) extends the classic Markov Decision Process (MDP) by addressing situations where the state of the system is only *partially observable* (Kaelbling et al. 1998). In real-world applications, decision makers often lack perfect information about the current state of the system, which complicates both control and information

gathering tasks. This partial observability introduces an inherent uncertainty that must be effectively managed for optimal decision making. The POMDP framework is uniquely suited to this challenge, as it allows decisions to be made based on probabilistic estimates of the state of the system, known as *belief states*.

The significance of solving POMDPs lies in their applicability to a wide range of complex real-world problems. POMDPs have been applied in fields such as reliability and maintenance (Cassandra 1998, Papakonstantinou and Shinozuka 2014, Kim et al. 2018, Song et al. 2022), aircraft collision avoidance (Temizer et al. 2010, Bai et al. 2012, Mueller and Kochenderfer 2016), and medical decision making (Vozikis and Goulionis 2009, Ayer et al. 2012, Zois 2016, Zhang and Wang 2022). In each of these domains, the ability to account for uncertain observations and dynamically update belief states based on new information is critical for success.

Although POMDPs offer a powerful theoretical framework, solving them exactly for large-scale problems remains computationally infeasible. This limitation arises from the *curse of dimensionality*, as the size of the reachable belief space grows exponentially with the number of states and observations (Papadimitriou and Tsitsiklis 1987, Madani et al. 1999). As a result, approximate methods have been developed to provide practical solutions. Among these, one of the most successful is the Point-Based Value Iteration (PBVI) algorithm (Pineau et al. 2003), which approximates the value function by focusing on a representative set of belief points rather than computing it over the entire belief space, which is a probability simplex, where each belief \mathbf{b} is a convex combination of the entire unit vectors corresponding to each state.

The purpose of this paper is to advance the current body of knowledge on solving finite-horizon POMDPs efficiently. In particular, we propose a novel method that accelerates the calculation of upper bounds in PBVI by using a Gaussian Process Upper Confidence Bound (GP-UCB). This method offers a data-efficient approximation strategy to improve the computational tractability of solving POMDPs in complex, large-scale environments. Our contribution is particularly significant in finite-horizon problems, where the dynamically changing value function requires repeated updates to both lower and upper bounds, the latter being computationally prohibitive to estimate. To this end, we employ Gaussian Process Regression (GPR) to learn a model of the upper bound estimates for a subset of the most informative belief points. The trained GPR model is then used to predict the entire upper bound convex hull, reducing the computational complexity while maintaining high-quality approximations. GP-UCB is proposed as a conservative estimate of the upper bound convex hull. The main contributions of our work are as follows.

1. The novel use of the GPR for the upper bound approximation, leveraging the non-parametric approximation capabilities of Gaussian Processes (GP)(Section 2).
2. The novel use of active learning for the iterative selection of the most informative belief points, referred to as support beliefs, to train the GPR. Training the GPR on this subset of the belief set reduces the computational cost from $\mathcal{O}(N^3)$ to $\mathcal{O}(Nd^2)$, where N is the number of belief points and $d \ll N$ is the number of support beliefs.

3. Theoretical guarantees, backed by formal proofs, showing that the upper bounds generated by our method converge to the true upper bound convex hull, leading to increasingly accurate approximations as the belief set expands.

Through extensive numerical experiments, we demonstrate that our method outperforms existing PBVI algorithms in both convergence speed and scalability across a range of complex problem domains and horizon lengths. Our experiments are conducted on well-established test problems, including instances with up to 90 states, 29 actions, and 3 observations. GP-UCB significantly improves computational efficiency in these settings, reducing computation time by 30%-60% on smaller problems and up to 99.7% on larger ones while maintaining the same gaps as the pure sawtooth projection method.

While POMDP formulations can involve even larger state and action spaces, our experiments exceed the scale of many Operations Research applications (e.g., in medical decision-making and maintenance planning), including recent ones, which typically consider formulations with at most 6 states, 3 actions, and 4 observations (Lin et al. 2004, Ayer et al. 2012, Liu et al. 2022, Hajjar and Alagoz 2023, Gong and Liu 2023, Li et al. 2023, Deep et al. 2023). By demonstrating strong performance on significantly larger problems than those addressed in these domains, our approach contributes to the OR literature by providing a scalable and computationally efficient tool for solving practical finite-horizon POMDPs.

Section 2 provides a description of the main problem that we consider, followed by a detailed explanation of the use of GPR to infer upper bounds in Section 3. Some readers may find the background given in Appendix A useful before proceeding with these sections since it provides an overview of POMDP solution approaches to date, with a focus on key components of state-of-the-art PBVI algorithms. We demonstrate our effectiveness claims in Section 5 through the use of an extensive set of numerical experiments on finite-horizon problems described in Section 4. The paper concludes with a summary and comments for future research directions in Section 6.

2 Problem Description

A POMDP is expressed by a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{O}, \Theta, \Omega, R, \mathbf{b}_0)$, where an agent occupies one of the possible states of the system $s \in \mathcal{S}$, which cannot be observed directly. The system changes from one state to the next after the agent takes an action $a \in \mathcal{A}$. Function $\Theta : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ represents the stochastic state transitions. Specifically, $\Theta(s, a, s') = P(s'|s, a)$ denotes the probability transition function of state changes from $s \in \mathcal{S}$ to $s' \in \mathcal{S}$ after performing the corresponding action $a \in \mathcal{A}$. Since the states are not observable directly, the agent makes certain observations $o \in \mathcal{O}$, which are imperfect projections of the states. The observation probability is defined by the function $\Omega : \mathcal{A} \times \mathcal{S} \times \mathcal{O} \rightarrow [0, 1]$. Specifically, $\Omega(a, s', o) = P(o|a, s')$ is the probability of observing $o \in \mathcal{O}$ after taking action $a \in \mathcal{A}$, under the true state $s' \in \mathcal{S}$. In addition, the agent’s action results in rewards. Function $R \in \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ denotes the reward function. Reward $R(s, a)$ is received after taking action $a \in \mathcal{A}$ in state $s \in \mathcal{S}$. The planning horizon T specifies the finite time steps during which the

agent seeks to maximize its cumulative reward. Let vector $\mathbf{b} = (b(1), b(2), \dots, b(|\mathcal{S}|))$ denote the belief state where $b(s)$ is the probability that the true system state is $s \in \mathcal{S}$. Starting from belief \mathbf{b} at time t , the agent updates belief to \mathbf{b}' at time $t + 1$ after executing action a and observing o as

$$b'(s') = \frac{P(o|a, s')}{P(o|\mathbf{b}, a)} \sum_{s \in \mathcal{S}} P(s'|s, a)b(s) = \frac{P(o|a, s') \sum_{s \in \mathcal{S}} P(s'|s, a)b(s)}{\sum_{s' \in \mathcal{S}} P(o|a, s') \sum_{s \in \mathcal{S}} P(s'|s, a)b(s)}, \forall s \in \mathcal{S}. \quad (2.1)$$

Naturally, $\sum_{s \in \mathcal{S}} b(s) = 1$. The initial belief, \mathbf{b}_0 , provides the probability distribution of the state at the beginning of the planning horizon. The beliefs derived from an initial belief, via a feasible sequence of actions and observations are called *reachable belief points*.

Figure 1 shows the belief states that are *reachable* from the initial belief state, $\mathbf{b}_0 = (0.5, 0.5)$ in one stage for the well-known *tiger problem* presented by Kaelbling et al. (1998), considering three actions, Listen (a_1), Open Left Door (a_2), Open Right Door (a_3), and two observations resulting from each action, under the assumption that opening a (left or right) door restarts the problem and resets the belief to $(0.5, 0.5)$. Appendix B contains a detailed description of the tiger problem.

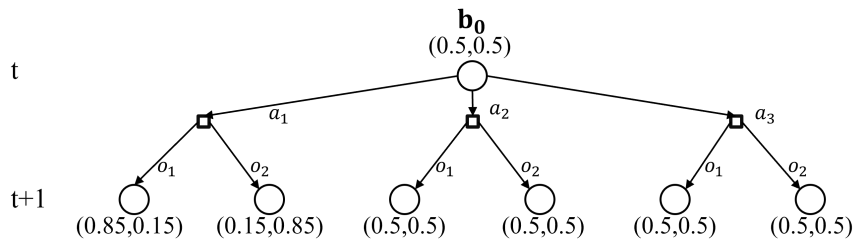


Figure 1: Tree structure of a two-stage tiger problem (Kaelbling et al. 1998). Circles represent belief states observed right before taking an action in each stage.

A policy π provides the sequence of actions to be taken over the planning horizon as a function of the belief state. The value function represents the expected cumulative reward obtained under the optimal policy, π^* . For a finite horizon problem, the following backward recursion computes the value function, $V_t(\mathbf{b})$, as,

$$V_t(\mathbf{b}) = \max_{a \in \mathcal{A}} \left[\sum_{s \in \mathcal{S}} b(s)R(s, a) + \sum_{o \in \mathcal{O}} P(o|\mathbf{b}, a)V_{t+1}^{\pi^*}(\mathbf{b}') \right], \forall \mathbf{b} \in \mathcal{B}_{\sqcup}, \sqcup \in \{t, \infty, \dots, \mathcal{T} - \infty\}, \quad (2.2)$$

where \mathcal{B}_{\sqcup} denotes the $|\mathcal{S}|$ -dimensional belief space. The optimal policy for belief state \mathbf{b} at time period t is defined as

$$\pi_t^*(\mathbf{b}) = \arg \max_{a \in \mathcal{A}} \left[\sum_{s \in \mathcal{S}} b(s)R(s, a) + \sum_{o \in \mathcal{O}} P(o|\mathbf{b}, a)V_{t+1}(\mathbf{b}') \right]. \quad (2.3)$$

Sondik (1971) showed that the value function is represented exactly by a piecewise-linear and convex function such that the value function for a specific period is represented by a set of $|\mathcal{S}|$ -

dimensional α -vectors, $\Gamma_t = \{\alpha^0, \alpha^1, \dots\}$. That is,

$$V_t(\mathbf{b}) = \max_{\alpha \in \Gamma_t} \sum_{s \in \mathcal{S}} b(s) \alpha(s), \quad t \in \{0, 1, \dots, T-1\}. \quad (2.4)$$

Hence, the value function given in Eqn. (2.2) can be rewritten as

$$V_t(\mathbf{b}) = \max_{a \in \mathcal{A}} \left[\sum_{s \in \mathcal{S}} b(s) R(s, a) + \sum_{o \in \mathcal{O}} P(o|\mathbf{b}, a) \max_{\alpha_{t+1} \in \Gamma_{t+1}} \left(\sum_{s' \in \mathcal{S}} b'(s') \alpha_{t+1}(s') \right) \right]. \quad (2.5)$$

This representation allows the value function to be computed recursively using a set of α -vectors. At each step, the *backup operation* updates the set of α -vectors, Γ_t based on the current belief state \mathbf{b} . Specifically, Γ_t is updated with

$$\Gamma_t = \bigcup_{\mathbf{b} \in \mathcal{B}_t} \text{Backup}(\mathbf{b}, t), \quad (2.6)$$

where the backup operation for each belief \mathbf{b} is given by

$$\text{Backup}(\mathbf{b}, t) = \arg \max_{g_a^b, \forall a \in \mathcal{A}} [\mathbf{b} \cdot g_a^b], \quad (2.7)$$

where g_a^b for action a is computed as

$$g_a^b(s) = \begin{cases} R(s, a) + \sum_{o \in \mathcal{O}} \arg \max_{\alpha_{t+1} \in \Gamma_{t+1}} \sum_{s \in \mathcal{S}} b(s) \sum_{s' \in \mathcal{S}} P(s'|s, a) P(o|a, s') \alpha_{t+1}(s'), & t < T, \\ R(s, a), & t = T. \end{cases} \quad (2.8)$$

PBVI Algorithm: A key idea of PBVI algorithms is sampling a representative set of belief points to approximate the value function. For these sampled belief points, both a lower bound \underline{V}_t and an upper bound \overline{V}_t of the value function are maintained. The gap between the lower and upper bounds at reachable belief points affects the value function approximation. Thus, the PBVI algorithm implemented by Walraven and Spaan (2019) selects the belief points with the largest gap, focusing on areas that can most effectively improve the value function approximation. The approach prioritizes the sampling of belief points where the gap is maximum in the time step $t+1$.

Algorithm 1 outlines the PBVI algorithm incorporating lower and upper bound calculations. Following Walraven and Spaan (2019), the initial belief set includes corner beliefs, which are the belief points where all probability mass is concentrated on a single state. Mathematically, these are the unit vectors \mathbf{w}_s , where \mathbf{w}_s is a vector with 1 in the s^{th} entry and 0 elsewhere. These points define the boundaries of the belief space, and hence are useful for initializing the PBVI and defining the reachable belief set.

We now discuss the computation of $\underline{V}_t(\mathbf{b})$ and $\overline{V}_t(\mathbf{b})$ for a belief point, $\mathbf{b} \in \mathcal{B}_\square$. While the former can be easily obtained by taking the best α -vector at any belief point (Hauskrecht 2000), obtaining

Algorithm 1 Point-based Value Iteration (PBVI) Algorithm, generalized from Walraven and Spaan (2019), Spaan and Vlassis (2005), Lovejoy (1991a)

```

1: Input: POMDP model, initial belief point  $\mathbf{b}_0$ , initial belief set  $\mathcal{B}_\sqcup$  containing all corner beliefs for
    $t \in \{0, 1, \dots, T-1\}$ , and convergence threshold,  $\epsilon$ 
2: Output: Lower bound  $\underline{V}_t(\mathbf{b})$  and upper bound  $\overline{V}_t(\mathbf{b})$  of value function for  $\forall \mathbf{b} \in \mathcal{B}_\sqcup$  for  $t \in \{0, 1, \dots, T-1\}$ 
3: Initialize  $\underline{V}_t(\mathbf{b})$  and  $\overline{V}_t(\mathbf{b})$  for  $\forall \mathbf{b} \in \mathcal{B}_\sqcup$  for  $t \in \{0, 1, \dots, T-1\}$  using Eqn. (2.5) and Eqn. (C.3)
4: while  $\overline{V}_0(\mathbf{b}_0) - \underline{V}_0(\mathbf{b}_0) > \epsilon$  do
5:   Choose a sampling method: Max-Gap, Random or Fixed-Grid
6:   if Sampling method = “Max-Gap” then
7:     Set  $\mathbf{b} = \mathbf{b}_0$ 
8:     for  $t = 0$  to  $T - 2$  do
9:       Choose the action  $a \leftarrow \arg \max_{a \in \mathcal{A}} (R(a) \cdot \mathbf{b} + \sum_{o \in \mathcal{O}} P(o|\mathbf{b}, a) \cdot \overline{V}_{t+1}(\mathbf{b}'))$ 
10:      Choose the observation  $o \leftarrow \arg \max_{o \in \mathcal{O}} (\overline{V}_{t+1}(\mathbf{b}) - \underline{V}_{t+1}(\mathbf{b}))$ 
11:      Sample new belief,  $\mathbf{b}_{\text{new}}$  using the selected  $a, o$  with Eqn. (2.1) and add to  $\mathcal{B}_{\sqcup+\infty}$ 
12:      Set  $\mathbf{b} = \mathbf{b}_{\text{new}}$ 
13:    end for (The max-gap method from Walraven and Spaan (2019))
14:   else if Sampling method = “Random” then
15:     for  $t = 1$  to  $T - 1$  do
16:       Sample new belief,  $\mathbf{b}_{\text{new}}$ , uniformly from belief space and add to  $\mathcal{B}_\sqcup$ , e.g., random sampling
       method proposed by Spaan and Vlassis (2005)
17:     end for
18:   else if Sampling method = “Fixed-Grid” then
19:     Use fixed belief grid for  $\mathcal{B}_\sqcup$  for  $t \in \{0, 1, \dots, T-1\}$ , e.g., fixed grid proposed by Lovejoy (1991a)
20:   end if
21:   for  $t = T - 1$  to  $0$  do
22:     Update  $\underline{V}_t(\mathbf{b})$  for  $\forall \mathbf{b} \in \mathcal{B}_\sqcup$  using Eqn. (2.5)
23:     Update  $\overline{V}_t(\mathbf{b})$  for  $\forall \mathbf{b} \in \mathcal{B}_\sqcup$  with Eqn. (C.3) using an approximation method, e.g., sawtooth
24:     Prune dominated  $\alpha$ -vectors
25:   end for
26: end while

```

the upper bound is computationally more challenging. In finite-horizon POMDPs, the challenge is even greater because the value function evolves dynamically over time (Smallwood and Sondik 1973). Unlike in infinite-horizon settings—where a stationary policy optimizes the value function indefinitely—finite-horizon problems require recomputing the value function at each time step, as the number of remaining decisions directly influences both immediate and future rewards (Pineau et al. 2003). This makes the calculation of upper bounds particularly expensive, since upper bounds must be recomputed at each time step to reflect the evolving decision horizon. Determining an exact upper bound requires optimistically evaluating all possible future outcomes, significantly increasing computational cost (Walraven and Spaan 2019, Smith and Simmons 2005).

Lovejoy (1991a) presented some of the earliest implementation of upper bounds using linear interpolation method using a grid of belief points obtained via Freudenthal triangulation. Subsequently, Hauskrecht (2000) presented an upper bound based on the convex hull projection while considering the upper bound improvement with incremental addition of new belief points. Since the value function is piecewise convex, the upper bound for a new belief point is obtained by projecting the belief point to the convex hull obtained by the existing belief-upper bound value pairs. For the

remainder of the paper, we refer to this as *convex hull*. Linear programming is used to identify the best convex hull by minimizing the linear combination of the existing upper bounds. Given the repetitive nature of updating upper bounds, the computational cost of executing a large number of linear programming steps quickly becomes intractable.

Hauskrecht (2000) proposed an interpolation approach for efficient approximation of the convex hull, referred to as *sawtooth*, to reduce computational complexity. Given an arbitrary belief $\mathbf{b} \in \mathcal{B}$ and corner beliefs, the sawtooth projection $\bar{v}(\cdot)$ provides the upper bound approximation for any new belief point. The detailed description of the sawtooth projection algorithm can be found in Appendix C.

To further demonstrate the sawtooth projection, consider the tiger problem again. In each subfigure of Figure 2, the orange dots represent belief points in the current belief set \mathcal{B} with known upper bounds. In the tiger problem, corner beliefs correspond to states where the tiger is surely behind either the left or the right door (i.e., belief states $(1.00, 0.00)$ or $(0.00, 1.00)$). Connecting non-corner orange belief points with these corner beliefs forms downward-pointing triangles, illustrated by the orange dashed lines. When a new belief \mathbf{b}' outside \mathcal{B} is encountered, its $\bar{v}(\mathbf{b}')$ is estimated by projecting it onto these connecting lines, selecting the projection with the smallest value (marked with red dots). For example, when updating the upper bound for belief points at $t = 3$, such as $(0.5, 0.5)$, Eqn. (C.3) is applied, requiring $\bar{v}(\cdot)$ for reachable beliefs at the subsequent stage, $t = 4$. In this case, the reachable belief points are $(0.5, 0.5)$, $(0.15, 0.85)$, and $(0.85, 0.15)$. Figure 2(a) shows that $(0.15, 0.85)$ is not part of the current belief set, so $\bar{v}[(0.15, 0.85)]$ is approximated using a sawtooth projection. By projecting $(0.15, 0.85)$ onto the dashed orange lines, the smallest projection value is selected, shown as the red point in Figure 2(a). For other reachable belief points whose upper bounds are known (orange points), the values are used as $\bar{v}(\cdot)$ directly. This process repeats as we move back to $t = 2$, $t = 1$, and $t = 0$, applying sawtooth projections to any reachable points with unknown upper bounds, as depicted by the red points in Figures 2(b)-(d).

The sawtooth projection method offers a trade-off between computational efficiency and accuracy. Using geometric properties, it efficiently approximates the convex hull, avoiding the extensive recalculations required by heuristic-based approaches such as HSVI (Smith and Simmons 2005) and SARSOP (Kurniawati et al. 2008). However, the sawtooth projection still presents significant computational complexity due to the repeated updates required to refine the upper bound each time a new belief point is added to the belief set. Initially, expanding the belief set results in significant improvements in the upper bound. But as the belief set grows, the improvement diminishes, although the number of updates, and consequently the computational complexity, continues to increase (Hauskrecht 2000).

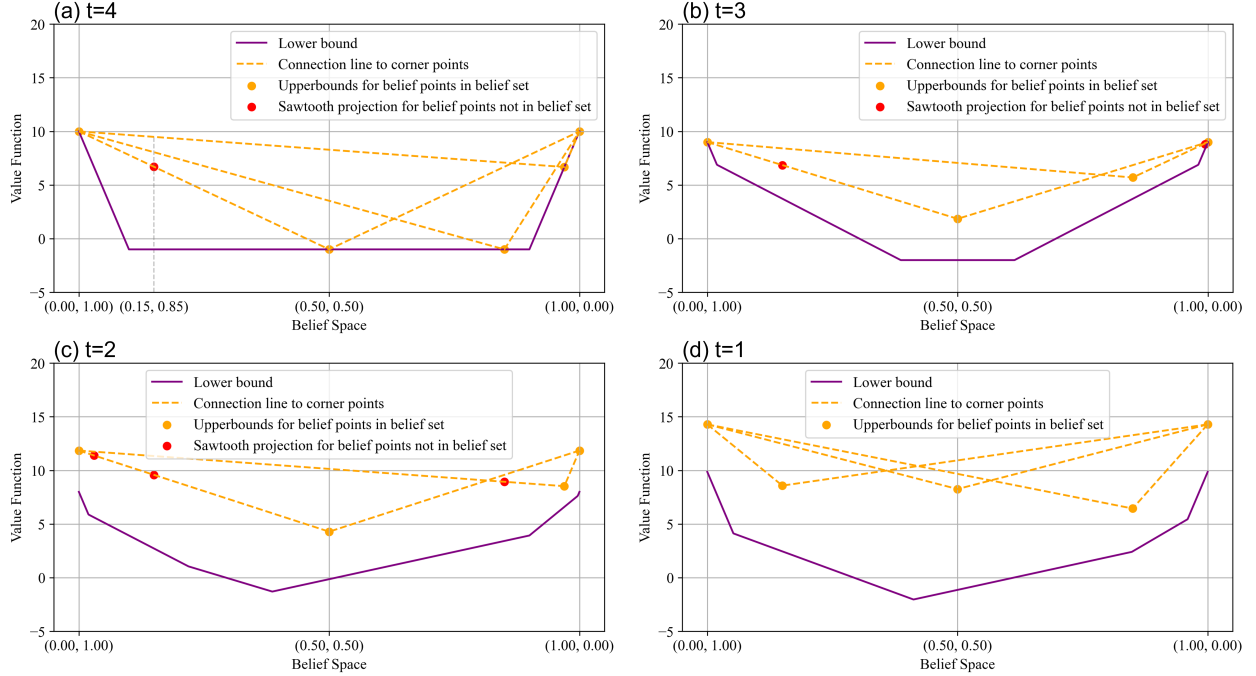


Figure 2: Approximation of the upper bound using sawtooth projection across different time stages at iteration 2 in the tiger problem.

3 Upper Bound Prediction using Gaussian Process Regression

Here, we introduce an alternative approach that approximates the upper bound convex hull by fitting a GPR to the sawtooth projections of only a subset of belief points that we refer to as the *support beliefs*. As noted in Hauskrecht (2000), the change in upper bounds as the belief set expands is only marginal. As a result, an “informative” subset of the belief set is typically sufficient to model the convex hull.

Recall that under the PBVI algorithm presented by Walraven and Spaan (2019), at any stage of upper bound updates for a set of belief points, we execute the backup operation as shown in Eqn. (C.3), which involves computing the sawtooth projection for all the reachable beliefs. Instead, we compute the sawtooth projection only for the support beliefs and predict for the rest of the belief points using a GPR. The support belief set is initialized with $|\mathcal{S}| + 1$ random belief points and iteratively expanded using an approximate linear dependence criterion discussed below. Finally, to ensure that our approximation provides a conservative estimate of the convex hull, we consider the upper confidence bound of the GPR during the inference stage.

GPR learns an interpolation of the convex hull using a set of m arbitrary reachable belief points and the corresponding sawtooth projections $\{\mathbf{b}_i, \bar{v}\}_{i=1}^m$. In this framework, the sawtooth projections $\bar{v}(\mathbf{b}_i)$ represent the noisy estimates of the convex hull, $h(\mathbf{b}_i)$, such that $\bar{v}(\mathbf{b}_i) = h(\mathbf{b}_i) + \epsilon_i$, where $\epsilon_i \sim \mathcal{N}(\mathbf{0}, \sigma_\epsilon^2)$ is assumed to be normally distributed noise with variance σ_ϵ^2 , representing the variability in the overestimation of the sawtooth projections. Hence, the sawtooth projection of

the training set follows a multivariate Gaussian distribution, i.e., $\bar{\mathbf{v}}_m \equiv [\bar{v}(\mathbf{b}_1), \dots, \bar{v}(\mathbf{b}_m)]^T \sim \mathcal{N}(\mathbf{0}, \mathbf{K}(\mathbf{B}_m, \mathbf{B}_m) + \sigma_\epsilon^2 \mathbf{I})$, where $\mathbf{B}_m = \{\mathbf{b}_i\}_{i=1}^m$ is the support belief set and $\mathbf{K}(\mathbf{B}_m, \mathbf{B}_m)$ denotes the covariance matrix. For any new belief point \mathbf{b}_k , the joint distribution of known upper bounds $\bar{\mathbf{v}}_m$ and estimated upper bound $\hat{v}(\mathbf{b}_k)$ is Gaussian, i.e.

$$\begin{bmatrix} \bar{\mathbf{v}}_m \\ \hat{v}(\mathbf{b}_k) \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K}(\mathbf{B}_m, \mathbf{B}_m) + \sigma_\epsilon^2 \mathbf{I} & \mathbf{K}(\mathbf{B}_m, \mathbf{b}_k) \\ \mathbf{K}(\mathbf{b}_k, \mathbf{B}_m) & \mathbf{K}(\mathbf{b}_k, \mathbf{b}_k) \end{bmatrix}\right) \quad (3.1)$$

where $\bar{\mathbf{v}}_m$ represents the upper bound at the current belief set \mathbf{B}_m , and $\hat{v}(\mathbf{b}_k)$ is the estimated upper bound at the new belief point \mathbf{b}_k . Conditioning on the joint distribution yields the predicted mean value $\mu(\mathbf{b}_k)$ of the convex hull at \mathbf{b}_k with the corresponding variance $\sigma^2(\mathbf{b}_k)$ at a new belief point \mathbf{b}_k as

$$\mu(\mathbf{b}_k) = \mathbf{K}(\mathbf{B}_m, \mathbf{b}_k)^T [\mathbf{K}(\mathbf{B}_m, \mathbf{B}_m) + \sigma_\epsilon^2 \mathbf{I}]^{-1} \bar{\mathbf{v}}_m, \quad (3.2)$$

$$\sigma^2(\mathbf{b}_k) = \mathbf{K}(\mathbf{b}_k, \mathbf{b}_k) - \mathbf{K}(\mathbf{B}_m, \mathbf{b}_k)^T [\mathbf{K}(\mathbf{B}_m, \mathbf{B}_m) + \sigma_\epsilon^2 \mathbf{I}]^{-1} \mathbf{K}(\mathbf{B}_m, \mathbf{b}_k). \quad (3.3)$$

Let us define $\boldsymbol{\beta} = [\mathbf{K}(\mathbf{B}_m, \mathbf{B}_m) + \sigma_\epsilon^2 \mathbf{I}]^{-1} \bar{\mathbf{v}}_m$ such that Eqn. (3.2) is rewritten as the following linear combination of kernel functions $k(\mathbf{b}_i, \mathbf{b}_k) = [\mathbf{K}(\mathbf{B}_m, \mathbf{b}_k)]_i$,

$$\mu(\mathbf{b}_k) = \sum_{i=1}^m \beta_i k(\mathbf{b}_i, \mathbf{b}_k) = \sum_{i=1}^m \beta_i \langle \phi(\mathbf{b}_i), \phi(\mathbf{b}_k) \rangle, \quad (3.4)$$

where ϕ is a mapping of the points in the belief space to Hilbert space. This form of the mean GPR predictor shows that it is, in fact, a linear predictor of the convex hull in the Hilbert space. As such, if an arbitrary belief point \mathbf{b}_k satisfies $\phi(\mathbf{b}_k) = \sum_{i=1}^m \beta_i \phi(\mathbf{b}_i)$, then the mean GPR prediction corresponding to \mathbf{b}_k could be inferred directly from the beliefs $\mathbf{b}_1, \dots, \mathbf{b}_m$. However, unless for the cases when the belief space is low dimensional or the belief $\mathbf{b}_k = \mathbf{b}_j, j \leq m$, $\phi(\mathbf{b}_k)$ will be linearly independent of $\{\phi(\mathbf{b}_i)\}_{i=1}^m$. Although strong linear dependency does not hold, Engel et al. (2004) showed that it is possible to have an approximate linear dependency. To test whether a new belief \mathbf{b}_k satisfies the approximately linear dependent or ALD criterion, we define

$$\delta \stackrel{\text{def}}{=} \min_{\mathbf{a}} \left\| \sum_{j=1}^m a_j \phi(\mathbf{b}_j) - \phi(\mathbf{b}_k) \right\|^2 = k(\mathbf{b}_k, \mathbf{b}_k) - \mathbf{K}(\mathbf{B}_m, \mathbf{b}_k)^\top \mathbf{a} \leq \nu, \quad (3.5)$$

where $\mathbf{a} \equiv [a_1, a_2, \dots, a_m]^\top = \mathbf{K}^{-1}(\mathbf{B}_m, \mathbf{B}_m) \mathbf{K}(\mathbf{B}_m, \mathbf{b}_k)$, and ν is a threshold that determines the strength of ALD. Here, $\mathbf{K}^{-1}(\mathbf{B}_m, \mathbf{B}_m)$ is the inverse of the kernel matrix for the belief set \mathbf{B}_m . If $\delta \leq \nu$, then \mathbf{b}_k is considered ALD on the support belief \mathbf{B}_m . In this case, the support belief set remains unchanged. However, if $\delta > \nu$, \mathbf{b}_k is not linearly dependent on \mathbf{B}_m , in which case the support belief \mathbf{B}_m is expanded to include the new belief point \mathbf{b}_k , and the corresponding kernel matrices are updated accordingly.

Note, however, that one needs to repeat the GPR fitting every time a new belief point \mathbf{b}_k is

added to the belief set. This is because as new belief points are added, the sawtooth projections \bar{v}_m of existing beliefs in the training set \mathbf{B}_m get tighter, leading to a more accurate upper bound approximation of the value function. However, these updated projections are not directly available and must be recomputed.

To reduce computational load, we update the GPR fit using partially revised sawtooth projections, which are then used to obtain the updated mean prediction in Eqn. (3.2), rather than retraining from scratch when upper bound changes are small. Specifically, during initial iterations, when upper bounds change significantly, the GPR model is refitted at each iteration to ensure accuracy. However, as the upper bounds stabilize in later iterations, a single belief point is randomly selected from \mathbf{B}_m in each iteration, and its sawtooth projection is recomputed to update the GPR fit in Eqn. (3.2). Full updates of all sawtooth projections \bar{v}_m become unnecessary because upper bound improvements tend to stabilize after a few iterations, as also noted by Hauskrecht (2000) and Smith and Simmons (2012).

Figure 3 provides a visualization of estimating upper bounds using GPR in the tiger problem example. Figure 3(a) shows the GPR model fit in iteration 1, and Figure 3(b) presents the fitting at iteration 5. At iteration 1, the support belief set includes three belief points with known upper bounds, depicted as orange dots, and three additional points selected based on the ALD criterion, with their respective sawtooth projections shown as green x markers. The GPR model is trained using these six belief points. The resulting convex hull inferred from the GPR model is shown in a solid blue line together with the 1σ -confidence bound. For comparison, the solid red line shows the upper bound approximation obtained from the sawtooth projection. As the belief set expands over subsequent iterations, the GPR model is updated with the new belief points using the ALD criterion. At iteration 5, the GPR model is trained by incorporating two additional reachable belief points alongside five existing beliefs with known upper bounds. The updated fitting of the GPR model reflecting the refined upper bound approximation, is depicted in Figure 3(b).

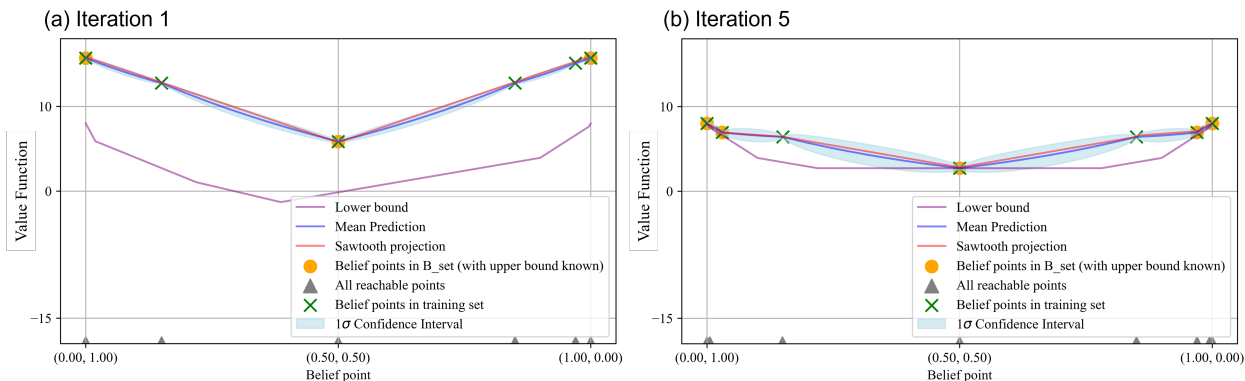


Figure 3: Upper bound estimation using GPR and sawtooth projection.

3.1 The Use of Gaussian Process Upper Confidence Bound

The last issue we address in this work is that of GP mean reversal and its impact on the convex hull prediction. The predicted mean tends to collapse toward the prior mean in regions with no or fewer training points. Since the prior mean is assumed to be 0 in the current case, the predicted convex hull would have a tendency to sag toward the zero line. The mean reversal effect is more pronounced with smoother kernel functions such as squared exponential and Matérn, but less so with exponential. An example comparison is shown in Figure 4.

Due to the mean reversal effect, the predicted convex hull may underestimate the true one. Hence, to overcome this limitation, we apply the Upper Confidence Bound (UCB) of the predicted convex hull as the upper bound $\bar{V}_t(\mathbf{b})$. For any belief point \mathbf{b} , the upper confidence bound $\bar{V}_{UCB}(\mathbf{b}) = \mu(\mathbf{b}) + \eta\sigma(\mathbf{b})$ provides a conservative estimate of the convex hull $h(\mathbf{b})$, where η is a confidence parameter. In the following, we theoretically show that the proposed upper bound $\bar{V}_{UCB}(\cdot)$ is a Probably Approximately Correct (PAC) estimate of the convex hull.

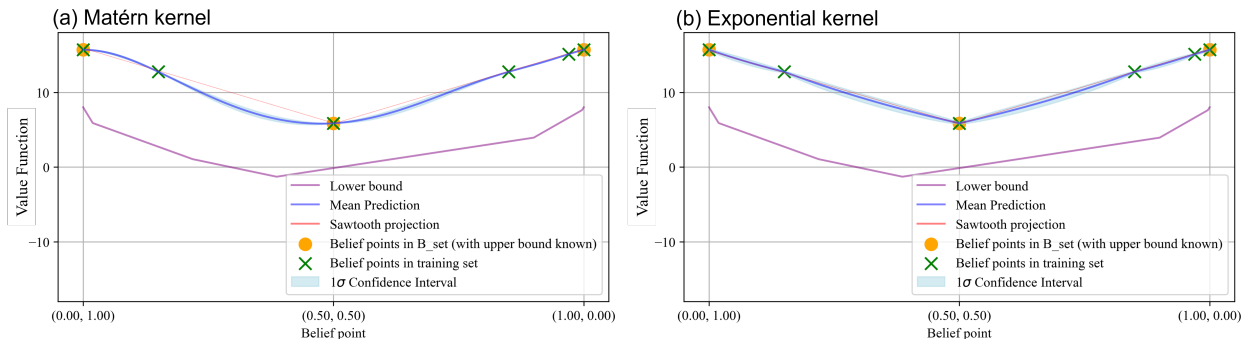


Figure 4: Upper bound estimation using Matérn and exponential kernels in the tiger problem iteration 1.

Theorem 1. Denoting the convex hull upper bound for belief point \mathbf{b} as $h(\mathbf{b})$, we have

$$P(|h(\mathbf{b}) - \mu(\mathbf{b})| \leq \eta\sigma(\mathbf{b})) \geq 1 - \delta \text{ for some } \delta \in (0, 1) .$$

In other words, as more belief points are sampled, $\sigma(\mathbf{b})$ reduces and the mean prediction $\mu(\mathbf{b})$ approaches the true convex hull $h(\mathbf{b})$.

Proof. Proof. Conditioned on a given set of belief points $\mathbf{B} = \{\mathbf{b}_i\}_{i=1}^m$ and their corresponding sawtooth projections for upper bounds $\bar{v}(\mathbf{b}_i)$, $i = 1, \dots, m$, the posterior distribution of the estimated upper bound using GPR is $\hat{v}(\mathbf{b}) \sim \mathcal{N}(\mu(\mathbf{b}), \sigma(\mathbf{b}))$, where $\mu(\mathbf{b})$ is the mean prediction and $\sigma(\mathbf{b})$ is the standard deviation.

As discussed in Hauskrecht (2000), the sawtooth projection $\bar{v}(\mathbf{b})$ overestimates the convex hull upper bound $h(\mathbf{b})$, ensuring $\bar{v}(\mathbf{b}) \geq h(\mathbf{b})$. Since the GPR model is trained on sawtooth projections, the posterior mean $\mu(\mathbf{b})$ progressively approximates $h(\mathbf{b})$ as additional belief points are sampled. Using the properties of Gaussian distribution, the normalized deviation of $h(\mathbf{b})$ from $\mu(\mathbf{b})$, relative

to the standard deviation $\sigma(\mathbf{b})$, satisfies

$$P\left(\frac{|h(\mathbf{b}) - \mu(\mathbf{b})|}{\sigma(\mathbf{b})} \leq \eta\right) \geq \frac{1}{\sqrt{2\pi}} e^{-\eta^2/2}.$$

Rearranging this inequality, we obtain $P(|h(\mathbf{b}) - \mu(\mathbf{b})| \leq \eta\sigma(\mathbf{b})) \geq 1 - \delta$, where $\delta = 1 - \frac{1}{\sqrt{2\pi}} e^{-\eta^2/2}$. □

Theorem 1 indicates that as more belief points are added, the uncertainty in the convex hull prediction, $\sigma(\mathbf{b})$ decreases, reflecting greater confidence in the GPR prediction. Consequently, the posterior mean $\mu(\mathbf{b})$ converges to the convex hull upper bound $h(\mathbf{b})$.

Corollary 1. *As $\mu(\mathbf{b})$ approaches $h(\mathbf{b})$, the proposed upper confidence bound $\bar{V}_{UCB}(\mathbf{b})$ becomes a conservative estimate of true upper bound, i.e.*

$$P(h(\mathbf{b}) \leq \bar{V}_{UCB}(\mathbf{b})) \geq \frac{1}{\sqrt{2\pi}} e^{-\eta^2/2}.$$

Proof. Proof. Consider the definition of the proposed upper confidence bound $\bar{V}_{UCB}(\mathbf{b})$: $\bar{V}_{UCB}(\mathbf{b}) = \mu(\mathbf{b}) + \eta\sigma(\mathbf{b})$. Theorem 1 ensures that $\mu(\mathbf{b})$ is close to $h(\mathbf{b})$ with high probability. Therefore, the probability that the upper bound convex hull $h(\mathbf{b})$ lies below the proposed upper confidence bound $\bar{V}_{UCB}(\mathbf{b})$ is $P(h(\mathbf{b}) \leq \bar{V}_{UCB}(\mathbf{b})) = P(h(\mathbf{b}) - \mu(\mathbf{b}) \leq \eta\sigma(\mathbf{b}))$. From Theorem 1, we know that $P(|h(\mathbf{b}) - \mu(\mathbf{b})| \leq \eta\sigma(\mathbf{b})) \geq 1 - \delta$. Substituting $\delta = 1 - \frac{1}{\sqrt{2\pi}} e^{-\eta^2/2}$, we get

$$P(h(\mathbf{b}) \leq \bar{V}_{UCB}(\mathbf{b})) \geq \frac{1}{\sqrt{2\pi}} e^{-\eta^2/2}.$$

Thus, as $\mu(\mathbf{b})$ converges to $h(\mathbf{b})$ with more belief points sampled and $\sigma(\mathbf{b})$ decreases, the proposed upper confidence bound becomes a probabilistically conservative estimate of the true upper bound $h(\mathbf{b})$. □

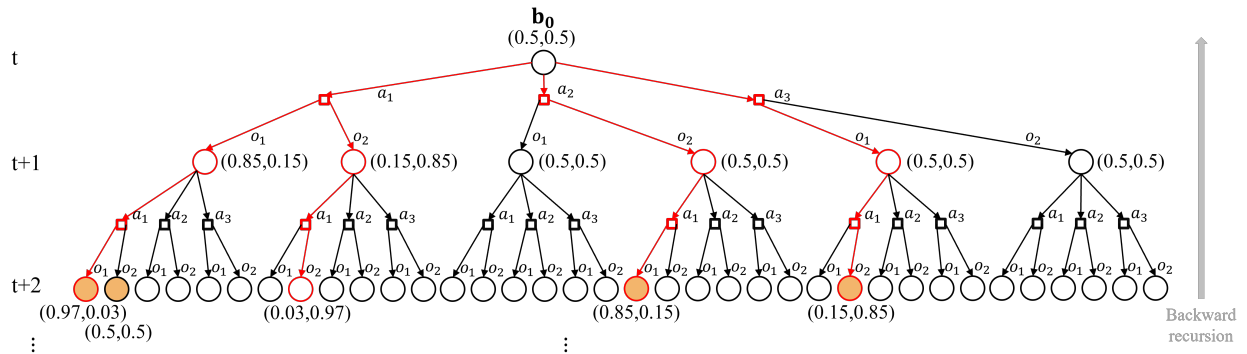


Figure 5: Estimation of upper bounds with the GPR model in the tiger problem.

Figure 5 shows how the GPR model estimates the upper bounds for the tiger problem. Starting from $(0.5, 0.5)$, the belief points reachable by $t + 2$ are depicted as circles, with sampled beliefs

shown on the red routes. The upper bounds of these sampled beliefs are updated using backward recursion. At $t + 2$, the GPR model is constructed using the most informative belief points in the belief set (shown with orange-shaded circles). In this example, one of the belief points at $t + 2$ (shown with a red edge but not orange-shaded) is not selected to train the GPR due to its lack of information gain. In addition, note that the belief point $(0.5, 0.5)$, which is not a belief point originally sampled, is selected to train the GPR model based on the ALD criterion, indicating it is an informative belief point. The trained GPR model is then used to estimate the upper bounds of other reachable belief points at $t + 2$ using the GP-UCB. Finally, we estimate upper bounds for belief points at $t + 1$ using the backward recursion given by Eqn. (C.3) in Appendix B.

Algorithm 2 outlines the steps for approximating an upper bound to the value function using our proposed GP-UCB approach.

Algorithm 2 GP-UCB Algorithm

```

1: Input: Planning horizon  $T$ , initial belief  $\mathbf{b}_0$ , threshold  $\epsilon$ 
2: Output:  $\overline{V}_0(\mathbf{b}_0)$ ,  $\underline{V}_0(\mathbf{b}_0)$ , and gap for initial belief  $\mathbf{b}_0$ 
3: Initialization: For  $t \in \{0, 1, \dots, T - 1\}$ , initialize  $\mathcal{B}_\square$  to include all corner points and  $\mathbf{b}_0$ , initialize training set  $\mathbf{B}_t$  to include all points in  $\mathcal{B}_\square$ , and initialize  $\Gamma_t$  using Eqn. (2.6).
4: while running time < time limit and gap >  $\epsilon$  do
5:   Expand  $\mathcal{B}_\square$  for  $t \in \{0, 1, \dots, T - 1\}$  using an sampling method, e.g., max-gap (lines 6-13 in Algo. 1), random (lines 14-17 in Algo. 1), fixed-grid (lines 18-20 in Algo. 1)
6:   for  $t = T - 1$  to 0 do
7:     Update  $\Gamma_t$  for  $\mathcal{B}_\square$  using Eqn. (2.6)
8:     if first iteration then
9:       Initialize GPR model for  $t$  with beliefs in  $\mathbf{B}_t$ 
10:    else
11:      if initial iterations or gap change >  $100 \cdot \epsilon$  or periodic check iteration then
12:        Update all upper bounds for beliefs in  $\mathbf{B}_t$  using sawtooth projection
13:        Fit the GPR model with updated upper bounds of the beliefs in  $\mathbf{B}_t$ 
14:      else
15:        Randomly select a belief  $\mathbf{b}$  in  $\mathbf{B}_t$ 
16:        Compute the updated sawtooth projection for  $\mathbf{b}$ 
17:        Update the GPR fit using revised sawtooth projections
18:      end if
19:    end if
20:    if  $\mathcal{B}_\square$  expanded with new points = True then
21:      Predict the covariance for new reachable beliefs using the GPR model
22:      if Covariance  $\geq$  ALD threshold then
23:        Expand  $\mathbf{B}_t$  with the new reachable belief
24:        Update GPR model using the expanded beliefs in  $\mathbf{B}_t$ 
25:      end if
26:    end if
27:    Update  $\overline{V}_t(\mathbf{b})$  for  $\forall \mathbf{b} \in \mathcal{B}_\square$  using Eqn. (C.3) with GPR model predictions
28:  end for
29:  Compute  $\overline{V}_0(\mathbf{b}_0)$  and  $\underline{V}_0(\mathbf{b}_0)$  for belief  $\mathbf{b}_0$  and update gap =  $\overline{V}_0(\mathbf{b}_0) - \underline{V}_0(\mathbf{b}_0)$ 
30: end while
31: return  $\overline{V}_0(\mathbf{b}_0)$ ,  $\underline{V}_0(\mathbf{b}_0)$ , and gap for the  $\mathbf{b}_0$ 

```

4 Numerical Experiments

In this section, we present the experimental evaluation of our approach using five commonly used test problems from pomdp.org (Cassandra 2025). These examples are selected to compare the performance of the algorithm on a range of problem sized, including varying numbers of states, actions, observations, and horizon length.

Table 1 summarizes the environmental variables for the five test problems used in our experiments. Since we focus on the finite planning horizon, we consider undiscounted versions of the problems (i.e., discount factors are set to 1). Additionally, we set the initial belief points for all examples as the center point, i.e., $(1/|\mathcal{S}|, 1/|\mathcal{S}|, \dots, 1/|\mathcal{S}|)$, which indicates that all states are equally likely at $t = 0$. To ensure robustness and account for variability due to the random selection of belief point in updating the GPR model (lines 15-17 in Algorithm 2) or random sampling, each experiment is repeated 10 times. All experiments were implemented in Python and solved on a Dell Desktop XPS 8940 with an Intel Core i5-11400 Processor and 40.0 GB RAM to ensure a fair comparison.

Table 1: Problem size variables for the test problems

	$ \mathcal{S} $	$ \mathcal{A} $	$ \mathcal{O} $	T
ChengD51	5	3	3	10, 15, 20, 40
Network	7	4	2	10, 15, 20, 40
Query	27	3	3	10, 15, 20, 40
Hallway	60	5	21	10, 15, 20, 40
Aloha.30	90	29	3	10, 15, 20, 40

Experiments are terminated based on two stopping criteria: either exceeding 3000 seconds or reaching a target gap of g_a . The target gap g_a is dynamically determined after each iteration based on the magnitude of the upper bounds of the value function for the belief point \mathbf{b}_0 at stage $t = 0$. Specifically, g_a is computed as $g_a = \frac{\Lambda(\bar{V}_0(\mathbf{b}_0))}{10^\rho}$, where $\bar{V}_0(\mathbf{b}_0)$ is the upper bound of the value function for the initial belief point \mathbf{b}_0 at stage $t = 0$ after the latest iteration, and $\Lambda(\bar{V}_0(\mathbf{b}_0))$ rounds $\bar{V}_0(\mathbf{b}_0)$ up to the nearest power of 10. For example, if $\bar{V}_0(\mathbf{b}_0) = 281$, then $\Lambda(281) = 1000$, and if $\bar{V}_0(\mathbf{b}_0) = 64$, then $\Lambda(64) = 100$. This way, the precision parameter ρ directly controls the stringency of g_a . Larger values of ρ yield smaller gaps, resulting in more precise approximations. We set the precision level, ρ , as 5 in our experiments. We set the uncertainty level, η , of the upper confidence bound to 1, and the accuracy level, ν , of the ALD condition to 10^{-5} . We set the initial iteration phase in line 11 of Algorithm 2 at 5 iterations and perform periodic checks every 5 iterations in our experiments.

To compare and contrast the performance of the different algorithms, we focus on two primary metrics: (i) the lower bound of the value function for the initial belief point \mathbf{b}_0 at stage $t = 0$, and (ii) the gap between the upper and lower bounds for \mathbf{b}_0 . For experiments involving random selection processes, such as the random selection of belief points in GPR model updates in the GP-UCB approach, both metrics of lower bound and gap are averaged over 10 runs to account for

variability and to ensure reliable comparisons.

In our study, we evaluate the performance of our GP-UCB approach in comparison to the sawtooth projection method under three different belief point sampling strategies: (1) max-gap sampling based on the work from Walraven and Spaan (2019), (2) random sampling proposed by Spaan and Vlassis (2005), and (3) fixed-grid sampling from Lovejoy (1991a). These strategies are chosen to highlight different aspects of belief point selection, with max-gap focusing on regions of high uncertainty in the value function, random sampling providing a baseline of unbiased coverage, and fixed-grid offering a structured exploration of the belief space. Combining each sampling strategy with the two approximation methods results in a total of six experimental settings, enabling a comprehensive analysis of their relative performance under varying sampling approaches. Among these six settings, fixed-grid sampling is fully deterministic for both the GP-UCB and sawtooth approaches, while max-gap sampling is also deterministic when used with the sawtooth approach.

5 Performance Results

A key advantage of the GP-UCB approach is its ability to significantly reduce the computational time required for sawtooth projections. Figure 6 compares the growth in the number of sawtooth executions over iterations for the max-gap-sawtooth and max-gap-GPUCB approaches, for all test problems. The results indicate that although both methods show an increase in sawtooth executions as iterations progress, the growth is substantially slower for max-gap-GPUCB. Using GP-UCB to estimate the upper bounds for belief points minimizes the need for frequent and computationally expensive sawtooth calculations. This efficiency allows the GP-UCB approach to achieve smaller gaps within the same computational time, or reach the target gap significantly faster than max-gap-sawtooth.

Table 2 presents the experimental results for lower bounds and gaps at $h = 40$, based on the predefined stopping criteria. Full results for all problems under all tested planning horizons (i.e., 10, 15, 20, and 40 periods) are provided in Tables D.1 to D.5. Bolded numbers indicate the largest lower bounds (LB) and smallest gaps (Gap), highlighting the best-performing methods for these metrics. The values shown in parentheses represent the standard deviations reflecting the variability observed across multiple runs, where applicable. In addition to reporting the average gap, we also report on the worst-case gap observed across the ten runs to provide additional information on the robustness and consistency of the methods under varying conditions.

Upon comparing the results of GP-UCB and pure sawtooth projection under the max-gap belief expansion strategy, it is evident that GP-UCB consistently achieves superior or comparable performance. Specifically, in terms of the lower bound, GP-UCB reaches the same or higher values as the sawtooth method in all cases by the termination of the experiment runs. Looking at the gaps, GP-UCB outperforms the sawtooth method in all examples, with the exception of the network problem at $h = 10$ (full results are provided in Appendix D). In this particular case, GP-UCB requires only 6 additional seconds to reach the target gap, demonstrating no significant disadvantage. For

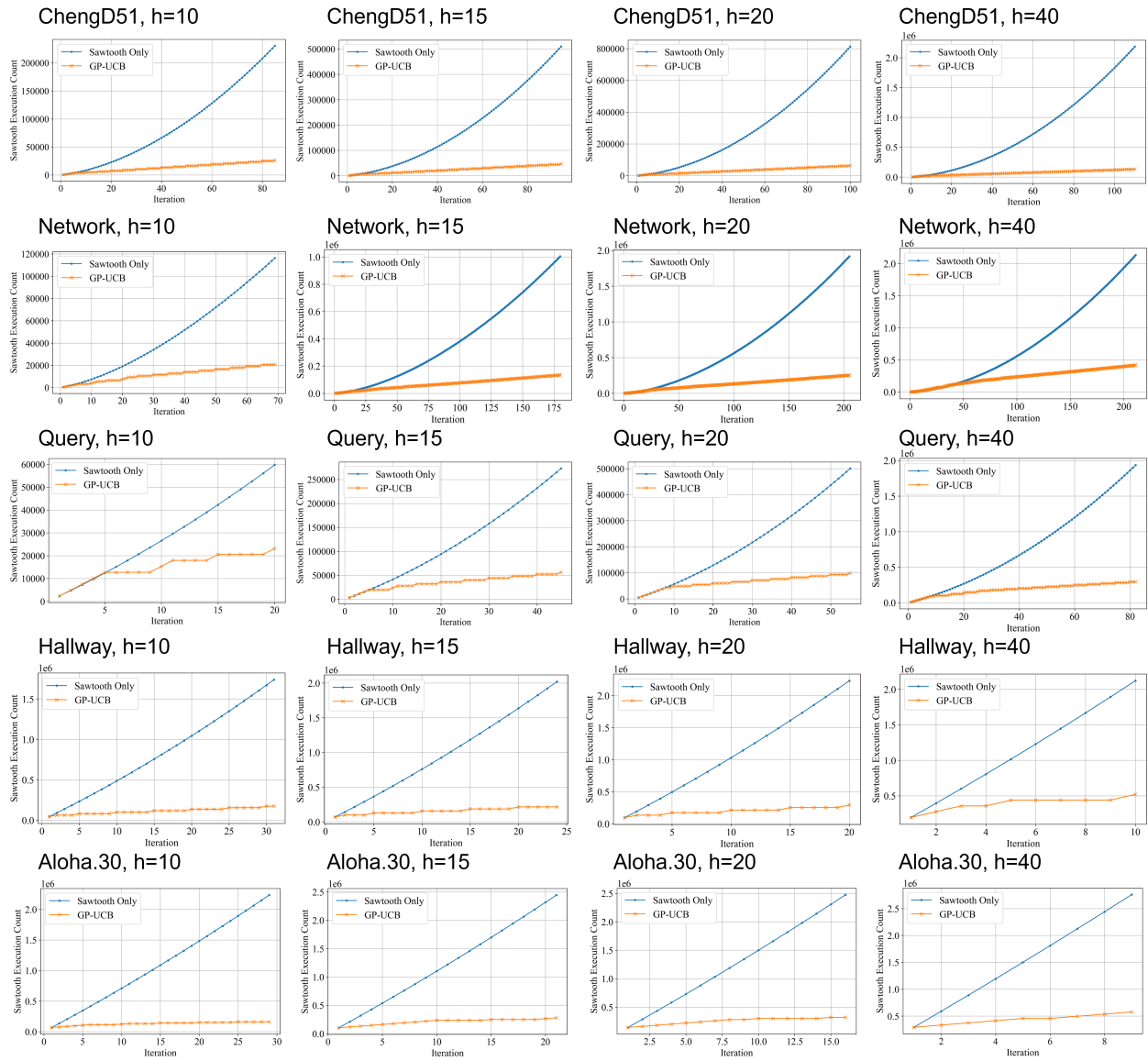


Figure 6: Comparison of the number of sawtooth executions over iterations between algorithms with pure sawtooth and with GP-UCB.

the remaining cases, GP-UCB either reaches the target gap faster or achieves a smaller gap when terminated at the 3000-second limit. Even when evaluating the worst-case gaps, the conclusions remain consistent with the observations for the average gaps across ten runs, further highlighting the efficiency of GP-UCB.

When evaluating the performance of sawtooth and GP-UCB under the random sampling strategy, we again observe consistent trends. For all test problems, GP-UCB achieves larger lower bounds and smaller gaps compared to the sawtooth method. This highlights the advantage of GP-UCB in improving the value function approximation even when belief points are sampled without a specific focus on regions of high uncertainty.

Finally, when comparing the GP-UCB and sawtooth projection methods under the fixed-grid belief sampling strategy, we observe that the fixed-grid strategy is unable to solve large-scale problems such as query, hallway, and aloha.30. This limitation arises due to the exponential growth in the number of grid points, leading to excessive computational and memory requirements. For the remaining two smaller problems, ChengD51 and network, the GP-UCB method consistently achieves the same lower bounds as the sawtooth method. However, GP-UCB demonstrates its superiority by achieving smaller gaps in all cases.

The above results demonstrate that the choice of belief expansion strategy significantly impacts performance. Among the three strategies evaluated (i.e., max-gap, random sampling, and fixed-grid sampling), fixed-grid strategy consistently unperformed and was unable to solve large-scale problems such as query, hallway, and aloha.30 due to its poor scalability to high-dimensional belief spaces. Max-gap sampling generally achieves the smallest gaps, emphasizing its effectiveness in reducing uncertainty in the belief space. However, for large-scale problems like aloha.30, random sampling achieves larger lower bounds than max-gap sampling, as it favors exploration of the belief space.

To further evaluate the efficiency of GP-UCB compared to the sawtooth projection method, we analyze the time required for GP-UCB to achieve the same (or smaller) gap as sawtooth at its terminating iteration under the max-gap belief expansion strategy. The results, summarized in Table 3, demonstrate significant time savings in most problems and planning horizons. For smaller problems like ChengD51 and network, GP-UCB achieves the sawtooth gaps in 30%-60% less time. For larger problems such as query, hallway, and aloha.30, the time savings are even more significant, reaching up to 99.7% for query and over 70% for aloha.30 at longer horizon lengths. While GP-UCB generally provides greater time savings at longer horizons, some variations exist across problems, such as Hallway. These differences likely stem from problem-specific structures and how efficiently GP-UCB generalizes across horizons. To capture these trends and exceptions, we present results for all planning horizons in Table 3.

To provide a visual representation of the experiment results, we plot the lower bounds and upper bounds of the value function during the execution of the max-gap-GP-UCB and max-gap-sawtooth methods at $h = 40$ in Figure 7. Full results are shown in Figure D.1. Consistent with the findings in Tables 2, the GP-UCB approach demonstrates strong performance in all test problems, achieving

Table 2: Experiment results for the five problems at $h = 40$. Bolded numbers indicate the largest upper bounds and smallest gaps, values in parentheses are standard deviations.

		Max-gap GP-UCB	Max-gap sawtooth	Random GP-UCB	Random sawtooth	Fixed-grid GP-UCB	Fixed-grid sawtooth
ChengD51	LB	261.713 ^(0.000)	261.713	261.713 ^(0.000)	261.713 ^(0.000)	261.711	261.711
	Gap	0.129 ^(0.002)	0.170	2.087 ^(0.058)	2.287 ^(0.092)	73.452	75.449
	Time (s)	3022.0 ^(9.3)	3028.8	3069.6 ^(26.9)	3019.3 ^(22.2)	2.3	0.5
	Worst Gap	0.131	0.170	2.170	2.447	73.452	75.449
Network	LB	592.274 ^(0.003)	592.274	591.878 ^(0.150)	591.719 ^(0.229)	582.704	582.704
	Gap	1.453 ^(0.018)	2.419	42.844 ^(1.844)	50.222 ^(1.503)	543.012	660.788
	Time (s)	3030.6 ^(36.2)	3020.2	3015.0 ^(15.7)	3020.2 ^(13.1)	3.2	1.5
	Worst Gap	1.488	2.419	45.825	53.214	543.012	660.788
Query	LB	120.327 ^(0.000)	120.327	120.327 ^(0.000)	120.327 ^(0.000)	NA	NA
	Gap	0.017 ^(0.001)	0.026	0.010 ^(0.001)	0.037 ^(0.000)	NA	NA
	Time (s)	3004.4 ^(5.6)	3049.5	3018.8 ^(32.9)	3027.4 ^(30.6)	NA	NA
	Worst Gap	0.018	0.026	0.010	0.037	NA	NA
Hallway	LB	1.930 ^(0.000)	1.912	1.612 ^(0.067)	1.295 ^(0.061)	NA	NA
	Gap	1.220 ^(0.000)	1.234	1.451 ^(0.070)	1.510 ^(0.061)	NA	NA
	Time (s)	2871.8 ^(18.9)	3030.7	3092.7 ^(66.6)	3069.2 ^(21.0)	NA	NA
	Worst Gap	1.220	1.234	1.571	1.638	NA	NA
Aloha.30	LB	282.129 ^(0.000)	282.039	283.338 ^(0.249)	282.329 ^(0.577)	NA	NA
	Gap	10.428 ^(0.003)	12.507	11.164 ^(0.632)	12.480 ^(0.932)	NA	NA
	Time (s)	2956.0 ^(27.2)	3086.8	3084.9 ^(45.8)	3163.7 ^(39.4)	NA	NA
	Worst Gap	10.435	12.507	12.113	14.570	NA	NA

Table 3: Time for GP-UCB to reach the same or smaller gap as sawtooth achieved at the terminating iteration for different problems using max-gap method for belief expansion

Problem	Planning Horizon	Sawtooth		GP-UCB		Time Reduction (%)
		Gap	Time (s)	Gap	Time (s)	
ChengD51	h=10	0.005	3022.4	0.005	1860.5	38.4%
	h=15	0.021	3010.9	0.021	1873.5	37.8%
	h=20	0.045	3032.9	0.044	1685.1	44.4%
	h=40	0.170	3028.8	0.167	1909.0	37.0%
Network	h=10	0.010	52.9	0.010	58.8	-11.2%
	h=15	0.039	3011.3	0.038	2060.9	31.6%
	h=20	0.346	3029.9	0.342	1183.4	60.9%
	h=40	2.419	3020.2	2.340	1369.1	54.7%
Query	h=10	0.001	80.7	0.001	5.7	92.9%
	h=15	0.002	3000.8	0.001	8.6	99.7%
	h=20	0.004	3025.6	0.003	11.1	99.6%
	h=40	0.026	3049.5	0.026	33.5	98.9%
Hallway	h=10	0.168	3034.5	0.166	2759.3	9.1%
	h=15	0.360	3082.8	0.358	1458.9	52.7%
	h=20	0.524	2883.4	0.524	2337.9	18.9%
	h=40	1.234	3030.7	1.222	2922.8	3.6%
Aloha.30	h=10	0.736	2962.7	0.735	2728.3	7.9%
	h=15	2.004	2912.5	1.978	776.9	73.3%
	h=20	4.003	3070.2	3.952	937.7	69.5%
	h=40	12.507	3086.8	12.419	740.2	76.0%

more favorable lower and upper bounds in most cases. For problems such as query and aloha.30, GP-UCB achieves significantly smaller upper bounds from the very beginning of the experiment runs, demonstrating its efficiency in approximating the value function in these cases. In contrast, for problems like the hallway problem, the reduction in upper bounds is less significant, indicating that certain problem structures may inherently limit the performance of GP-UCB.

One limitation of our approach is the need to select the kernel function. In these experiments, we used an exponential kernel for the GP-UCB approach. While this choice works well for most problems, larger problems like hallway and aloha.30 may benefit from alternative kernels such as the Matérn kernel, which can better generalize across larger belief spaces and avoid overfitting to the small initial data sets.

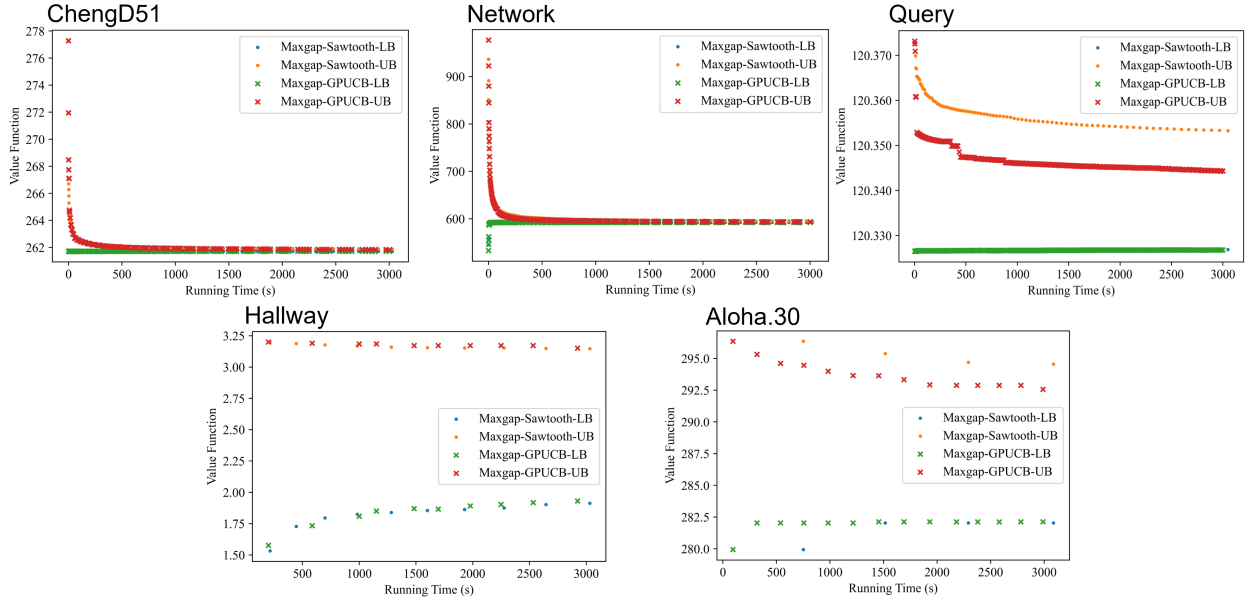


Figure 7: Upper and lower bounds of the value function for the five examples at $h = 40$. Dots represent upper and lower bounds from FiVI, while crosses indicate those from AUCB, color-coded as shown in the legend.

6 Conclusions

In this paper, we present GP-UCB, a novel approach to accelerate the upper bound estimation in point-based value iteration to solve finite-horizon POMDPs. By utilizing Gaussian Process Regression to approximate the upper bound convex hull across the belief space, our method reduces computational complexity without sacrificing solution quality. The approach selectively trains on the most informative subset of belief points, significantly enhancing solution efficiency and scalability. This innovation avoids redundant upper bound updates and provides theoretical guarantees for the convergence of the proposed GP-UCB to the true upper bound convex hull. Through extensive experiments on benchmark problems, we demonstrated that GP-UCB consistently accelerates

convergence and improves scalability, reducing computation time by 30%-60% on smaller problems and up to 99.7% on larger ones while maintaining the same upper bound gaps as the pure sawtooth projection method. These results underscore its potential as a powerful tool for solving practical large-scale POMDPs in finite-horizon settings.

In the proposed approach, the GPR model uses training data from sawtooth projections, which are computationally efficient for estimating upper bounds. However, a trade-off exists between computation time and data quality. Generating higher-quality data results in better approximations for upper bounds, but may take longer. Future research will focus on finding a way to balance this trade-off to improve GPR model performance. Another promising direction is optimizing the selection of kernels for GPR, as the choice of kernel directly impacts the quality of the approximation and computational efficiency. Investigating adaptive kernel selection strategies or using domain-specific kernels could enhance the flexibility and accuracy of the GPR model.

References

- Ayer T, Alagoz O, Stout NK (2012) OR forum—a POMDP approach to personalize mammography screening decisions. *Operations Research* 60(5):1019–1034.
- Bai H, Hsu D, Kochenderfer MJ, Lee WS (2012) Unmanned aircraft collision avoidance using continuous-state POMDPs. *Robotics: Science and Systems VII* 1:1–8.
- Cassandra AR (1998) A survey of POMDP applications. *Working notes of AAAI 1998 fall symposium on planning with partially observable Markov decision processes*, volume 1724.
- Cassandra AR (2025) Partially Observable Markov Decision Processes (POMDPs). <http://www.pomdp.org>, accessed: 04-Feb-2025.
- Cassandra AR, Littman ML, Zhang NL (1997) Incremental pruning: A simple, fast, exact method for partially observable Markov decision processes. *Proceedings of the Thirteenth Conference on Uncertainty in Artificial Intelligence (UAI)* 54–61.
- Cheng Hh (1988) Algorithms for partially observable markov decision processes. *Proceedings of the 24th IEEE Conference on Decision and Control*, 1209–1211 (IEEE).
- Deep A, Zhou S, Veeramani D, Chen Y (2023) Partially observable markov decision process-based optimal maintenance planning with time-dependent observations. *European Journal of Operational Research* 311(2):533–544.
- Engel Y, Mannor S, Meir R (2004) The kernel recursive least-squares algorithm. *IEEE Transactions on signal processing* 52(8):2275–2285.
- Gong J, Liu S (2023) Partially observable collaborative model for optimizing personalized treatment selection. *European Journal of Operational Research* 309(3):1409–1419.
- Hajjar A, Alagoz O (2023) Personalized disease screening decisions considering a chronic condition. *Management Science* 69(1):260–282.
- Hauskrecht M (2000) Value-function approximations for partially observable Markov decision processes. *Journal of artificial intelligence research* 13:33–94.
- Islam UJ, Paynabar K, Runger G, Iqbal AS (2024) Dynamic exploration–exploitation trade-off in active learning regression with bayesian hierarchical modeling. *IISE Transactions* 1–15.

- Kaelbling LP, Littman ML, Cassandra AR (1998) Planning and acting in partially observable stochastic domains. *Artificial intelligence* 101(1-2):99–134.
- Kim JW, Choi GB, Lee JM (2018) A POMDP framework for integrated scheduling of infrastructure maintenance and inspection. *Computers & Chemical Engineering* 112:239–252.
- Kurniawati H, Hsu D, Lee WS (2008) SARSOP: Efficient point-based POMDP planning by approximating optimally reachable belief spaces. *Robotics: Science and systems*, volume 2008 (Zurich, Switzerland).
- Li W, Denton BT, Morgan TM (2023) Optimizing active surveillance for prostate cancer using Partially Observable Markov Decision Processes. *European Journal of Operational Research* 305(1):386–399.
- Lin ZZ, Bean JC, White CC (2004) A hybrid genetic/optimization algorithm for finite-horizon, partially observed Markov decision processes. *INFORMS Journal on Computing* 16(1):27–38.
- Littman ML (1996) *Algorithms for sequential decision-making* (Brown University).
- Littman ML (2009) A tutorial on partially observable Markov decision processes. *Journal of Mathematical Psychology* 53(3):119–125.
- Littman ML, Cassandra AR, Kaelbling LP (1995) Learning policies for partially observable environments: Scaling up. *Proceedings of the Twelfth International Conference on Machine Learning (ICML)*, 362–370.
- Liu Z, Khojandi A, Li X, Mohammed A, Davis RL, Kamaleswaran R (2022) A machine learning-enabled partially observable markov decision process framework for early sepsis prediction. *INFORMS Journal on Computing* 34(4):2039–2057.
- Lovejoy WS (1991a) Computationally feasible bounds for partially observed Markov decision processes. *Operations Research* 39(1):162–175.
- Lovejoy WS (1991b) A survey of algorithmic methods for partially observed Markov decision processes. *Annals of Operations Research* 28(1):47–65.
- Madani O, Hanks S, Condon A (1999) On the undecidability of probabilistic planning and infinite-horizon partially observable Markov decision problems. *Aaai/iaai* 10(315149.315395).
- Monahan G (1982) A survey of partially observable markov decision processes: theory, models, and algorithms. *Management Science* 28(1):1–16.
- Mueller ER, Kochenderfer M (2016) Multi-rotor aircraft collision avoidance using partially observable Markov decision processes. *AIAA Modeling and Simulation Technologies Conference*, 3673.
- Papadimitriou CH, Tsitsiklis JN (1987) The complexity of Markov decision processes. *Mathematics of Operations Research* 12(3):441–450.
- Papakonstantinou KG, Shinozuka M (2014) Planning structural inspection and maintenance policies via dynamic programming and Markov processes. part ii: POMDP implementation. *Reliability Engineering & System Safety* 130:214–224.
- Pineau J, Gordon G, Thrun S (2003) Point-based value iteration: An anytime algorithm for POMDPs. *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, 1025–1030.
- Poupart P, Kim KE, Kim D (2011) Closing the gap: Improved bounds on optimal POMDP solutions. *Proceedings of the International Conference on Automated Planning and Scheduling*, volume 21.
- Shani G, Pineau J, Kaplow R (2013) A survey of point-based POMDP solvers. *Autonomous Agents and Multi-Agent Systems* 27(1):1–51.
- Smallwood RD, Sondik EJ (1973) The optimal control of partially observable markov processes over a finite horizon. *Operations Research* 21(5):1071–1088.

- Smith T, Simmons R (2004) Heuristic search value iteration for pomdps. *Proceedings of the 20th Conference on Uncertainty in Artificial Intelligence (UAI)*, 520–527 (AUAI Press).
- Smith T, Simmons R (2005) Point-based POMDP algorithms: Improved analysis and implementation. *Proceedings of the 21st Conference on Uncertainty in Artificial Intelligence (UAI)* 542–549.
- Smith T, Simmons R (2012) Point-based POMDP algorithms: Improved analysis and implementation. *arXiv preprint arXiv:1207.1412* .
- Sondik EJ (1971) *The optimal control of Partially Observable Markov Processes* (Stanford University).
- Sondik EJ (1978) The optimal control of Partially Observable Markov Processes over the infinite horizon: Discounted costs. *Operations Research* 26(2):282–304.
- Song C, Zhang C, Shafieezadeh A, Xiao R (2022) Value of information analysis in non-stationary stochastic decision environments: A reliability-assisted POMDP approach. *Reliability Engineering & System Safety* 217:108034.
- Spaan MT, Vlassis N (2005) Perseus: Randomized point-based value iteration for POMDPs. *Journal of artificial intelligence research* 24:195–220.
- Temizer S, Kochenderfer M, Kaelbling L, Lozano-Pérez T, Kuchar J (2010) Collision avoidance for unmanned aircraft using Markov decision processes. *AIAA Guidance Navigation and Control Conference*, 8040.
- Vozikis A, Goulionis JE (2009) Medical decision making for patients with parkinson disease under average cost criterion. *Australia and New Zealand Health Policy* 6(1).
- Walraven E, Spaan MT (2019) Point-based value iteration for finite-horizon POMDPs. *Journal of Artificial Intelligence Research* 65:307–341.
- Zhang NL, Liu W (1996) Planning in stochastic domains: Problem characteristics and approximation. Technical report, Technical Report HKUST-CS96-31, Department of Computer Science, Hong Kong University of Science and Technology.
- Zhang W, Wang H (2022) Diagnostic policies optimization for chronic diseases based on POMDP model. *Healthcare*, volume 10:2, 283 (MDPI).
- Zois DS (2016) Sequential decision-making in healthcare iot: Real-time health monitoring, treatments and interventions. *2016 IEEE 3rd World Forum on Internet of Things (WF-IoT)*, 24–29 (IEEE).

A Background on POMDP Problems and Solution Algorithms

Like MDPs, POMDPs identify control actions based on the current state of the system, but since the system state is only partially observable, decisions are made based on a *belief state*—a probability distribution over the possible system states that represents the current estimate of the system’s true state. This belief state is updated through systemic observations that provide probabilistic information about the system state at any given time. In addition to the usual MDP assumptions regarding action- and state-dependent transition probabilities, POMDP formulations assume a *likelihood function*, which gives the probability of observing a particular outcome for each system state. This allows for the computation of the next belief state as a function of the current belief state, the action taken, and the observation obtained (Littman 2009). Since the belief state is a continuous vector over the probability distribution of the system states, the belief space becomes continuous, making the solution of POMDPs much more difficult (Papadimitriou and Tsitsiklis 1987, Kaelbling et al. 1998, Madani et al. 1999).

The structural properties of the *value function* in POMDPs, for both finite- and infinite-horizon problems, are well understood. In both cases, the value function, which maps a belief state to the expected total reward, is *piecewise-linear and convex* in the belief space (Sondik 1971, Smallwood and Sondik 1973). This means that the value function can be represented as the maximum of a set of linear functions over the belief space, each associated with a specific action. These linear functions are known as *alpha-vectors*, where each alpha-vector corresponds to a specific action and encodes the expected future rewards for that action given the current belief state (Sondik 1971). For any belief state, the value is determined by selecting the alpha-vector that maximizes the expected reward. In the finite-horizon case, the set of alpha-vectors is updated stage by stage using backward induction to compute the value function at each time step (Smallwood and Sondik 1973). In infinite-horizon problems, the value function is stationary, and the alpha-vectors are iteratively refined until a fixed point is reached (Sondik 1978, Monahan 1982, Lovejoy 1991b). However, the number of belief states and corresponding alpha-vectors grows *exponentially* with the size of the state and observation spaces, leading to the well-known “curse of dimensionality.” This exponential growth makes the computation of exact solutions infeasible for large-scale problems, even when pruning techniques are employed to reduce the number of alpha-vectors (Sondik 1971, Monahan 1982, Cheng 1988, Littman 1996, Zhang and Liu 1996, Cassandra et al. 1997).

A.1 Point-Based Value Iteration

After the structural properties of POMDP value functions were established, and the challenges of solving them exactly for large problems became evident, researchers developed various approximate methods. Among the most prominent of these is the Point-Based Value Iteration (PBVI) algorithm, introduced by Pineau et al. (2003). PBVI was a major breakthrough because it significantly improved the scalability of POMDP solvers by focusing on a subset of belief points, rather than attempting to compute the value function over the entire continuous belief space.

The key insight behind PBVI is that not all belief states are equally important for decision making. Instead of considering the entire belief space, PBVI selects a representative set of belief points and performs computations on these points. By concentrating computational effort on a smaller set of critical belief states, PBVI drastically reduces the computational load while still providing a good approximation of the value function. This approach is particularly effective, as PBVI uses the structure of POMDPs where a limited set of belief points - especially those likely to be encountered under the optimal policy - can be sufficient to approximate the value function well enough for effective decision making (Pineau et al. 2003, Smith and Simmons 2005). This makes PBVI one of the most widely used algorithms to solve POMDPs, particularly for large and complex problems (Shani et al. 2013).

The main components of PBVI algorithms can be summarized as follows:

1. **Select Initial Belief Points:** A small set of representative belief points is initialized, typically using random sampling or simulations.
2. **Backup Operation:** The lower and upper bounds on the value function at each belief point are updated by considering the immediate reward and expected future rewards based on potential actions and observations.
3. **Expand Belief Set:** New belief points are added as needed to improve the approximation across the belief space.
4. **Convergence:** The process repeats until the difference between the upper and lower bounds falls below a set threshold, indicating convergence.

Since its introduction in 2003, various improvements have been made to each of the four key components of PBVI. For belief point selection, techniques such as heuristic sampling (Pineau et al. 2003) and more structured approaches such as reachability analysis (Kurniawati et al. 2008) have improved the initial set of belief points by focusing on those most relevant to the policy. The backup operation has seen optimizations with randomized belief backups, as in the Perseus algorithm (Spaan and Vlassis 2005), which selectively updates belief points to reduce computational cost. In terms of belief set expansion, advanced methods such as adaptive belief selection (Shani et al. 2013) and forward search techniques have improved the algorithm’s ability to refine the belief set efficiently.

For convergence, several approaches have been employed to accelerate the process. One such strategy is incremental pruning (Cassandra et al. 1997), which reduces the complexity of managing alpha-vectors by eliminating irrelevant vectors during each iteration. Another is heuristic search, which focuses computational resources on the most promising regions of the belief space, improving both convergence speed and efficiency (Smith and Simmons 2004). Additionally, bounded policy updates as implemented in Heuristic Search Value Iteration (HSVI) (Smith and Simmons 2005) ensure that the value function and policy are refined efficiently while maintaining convergence guarantees.

A.2 Upper and Lower Bounds

One of the most significant innovations for convergence and belief point selection is the introduction of upper and lower bounds on the value function. These bounds serve multiple purposes: they guide the selection of critical belief points by focusing on regions where the difference between the upper and lower bounds is largest, and they provide stopping criteria for the algorithm when the bounds converge sufficiently (Smith and Simmons 2005, Poupart et al. 2011). By integrating bounds into PBVI, convergence is not only accelerated, but the algorithm also gains the ability to focus more intelligently on belief points where improvement is needed most.

In PBVI algorithms, calculating the lower bounds is relatively straightforward and computationally inexpensive. For each belief point, the best *alpha-vector* (the one that yields the highest expected reward) is chosen, ensuring that the lower bound represents a conservative estimate of the value function. This method provides a reliable, minimal reward estimate without adding significant computational overhead. However, calculating upper bounds has traditionally been more complex and costly. Exact methods for the upper bound calculation involve projecting the belief point onto the convex hull of the belief/upper bound pairs. The tightest upper bound is obtained by minimizing the projection using linear programming methods. Although this provides the strongest possible upper bound, it is computationally expensive (Monahan 1982, Cheng 1988). The QMDP approach, which uses Q -values (i.e., the expected reward of taking an action in a state and following an optimal policy), simplifies the problem by assuming full observability after a single step, leading to an optimistic upper bound (Littman et al. 1995).

For such cases, approximate methods like those used in HSVI and Successive Approximations of the Reachable Space under Optimal Policies (SARSOP) algorithms offer more computational efficiency by deriving upper bounds from optimistic assumptions. In HSVI, the upper bound is updated by a heuristic search, exploring actions and observations optimistically, guiding the algorithm toward promising belief points (Smith and Simmons 2005). SARSOP focuses on reachable belief spaces to update upper bounds only for belief points likely to be encountered under an optimal policy, further reducing computational costs (Kurniawati et al. 2008). These heuristic-based methods are particularly useful in infinite-horizon POMDPs, where the value function is stationary, meaning upper bounds do not need frequent recalculation.

In finite-horizon POMDPs, the situation becomes more complex because the value function changes dynamically over time (Walraven and Spaan 2019). The upper bound must be recalculated at each time step, reflecting the evolving decision horizon. Calculating an exact upper bound for finite-horizon problems involves exploring all possible future outcomes optimistically, considering the best-case scenario at each stage, which leads to a high computational cost (Walraven and Spaan 2019, Smith and Simmons 2005). This process requires evaluating future rewards under all potential actions and observations for each belief point, further increasing the complexity.

Although PBVI has been effective for infinite-horizon problems, where a stationary policy optimizes the value function indefinitely, significant adjustments are necessary to apply it to finite-horizon problems. In infinite-horizon settings, the value function remains stationary, allowing algo-

rithms such as PBVI to refine the value function iteratively without accounting for the remaining time steps (Pineau et al. 2003). However, in finite-horizon problems, the value function changes since the number of remaining decisions affects both immediate and future rewards. As the horizon shortens, the value function increasingly reflects immediate rather than future rewards. This requires the algorithm to adapt at each step to account for the evolving optimal policy (Walraven and Spaan 2019). Efficient mechanisms, particularly for accurate upper and lower bounds, are crucial to ensuring that belief point selection and value function updates remain computationally feasible.

To address these challenges, Walraven and Spaan (2019) implemented efficient methods for handling upper and lower bounds in finite-horizon problems. By estimating the maximum possible value for each belief point and focusing on the widest gap between the upper and lower bounds, the algorithm prioritizes regions where the value function is the least accurate. This accelerates convergence and reduces computation time, while upper bounds serve as a stopping criterion, terminating the algorithm when convergence is reached, ensuring a near-optimal policy without excessive refinement.

A.3 Sawtooth Projection to Obtain Upper Bounds

Instead of using exact upper bounds, Hauskrecht (2000) introduced the so-called *sawtooth projections* to dynamically adjust the upper bounds at each step of the horizon. This approximation leverages the convex structure of the value function to reduce computational complexity while maintaining a reasonable level of approximation accuracy. Poupart et al. (2011) reported that sawtooth projections can significantly reduce computational costs while preserving a high degree of solution quality. Specifically, in the worst-case scenario, the computational complexity for the approximation decreases to $O(|\mathcal{B}||\mathcal{S}|)$, where \mathcal{B} is a sampled belief set in the $|\mathcal{S}|$ -dimensional space with $|\mathcal{B}| > |\mathcal{S}|$.

The effectiveness of sawtooth projection in finite-horizon settings was recently demonstrated by Walraven and Spaan (2019) in their Finite-horizon Value Iteration (FiVI) algorithm. FiVI outperforms both the original PBVI (Pineau et al. 2003) and the Heuristic Search Value Iteration (HSVI) (Smith and Simmons 2005) in terms of speed and scalability, especially in large finite-horizon problems. By focusing computational effort where it is most needed, the sawtooth projection allows FiVI to achieve faster convergence without sacrificing solution quality. Furthermore, their numerical results showed that FiVI scales better with both the horizon length and the complexity of belief spaces. This makes it a more versatile option for large-scale applications, and less sensitive to the size of state and action spaces (Walraven and Spaan 2019).

The primary computational complexity of the sawtooth projection stems from the need to repeatedly update the upper bound each time a new belief point is added to the belief set. As this set expands, the upper bound approximation of the value function improves, narrowing the gap with the lower bound. In the initial stages of belief expansion, adding new points significantly improves the upper bound. However, as more belief points are added, the marginal improvement becomes minimal, even though the number of updates—and thus the computational complexity—

increases (Hauskrecht 2000). This leads to upper bound updates that require computations that scale with the product of the cardinality of the belief set and the dimensionality of the belief state, resulting in a computational complexity of $O(|\mathcal{B}||\mathcal{S}|)$.

Hence, while sawtooth projection offers a computationally efficient alternative to computing the convex hull of upper bounds, its effectiveness remains limited, as the number of sawtooth projection operations grows rapidly with the expanding belief set. As a result, this approach is still insufficient to significantly reduce complexity and solve larger problems.

In this work, we model the convex hull of the upper bounds by fitting a subset of the belief points and their sawtooth projections using GPR. As such, for any new belief point, the upper bound interpolation is directly inferred from the posterior of the GPR predictions, precluding the need to solve any linear programs or sawtooth projection. Furthermore, at every iteration of the upper bound approximation, we fit the GPR only to a subset of the belief/upper bound pairs that are most informative, to reduce the computational cost of covariance matrix inversion.

Notably, multiple belief points lying on the same alpha-vector are redundant and do not provide additional information about the value function. To avoid such redundant belief points while training the GPR, we argue that the most informative belief points are those that are approximately linearly independent in a high-dimensional Hilbert space (Islam et al. 2024, Engel et al. 2004). Therefore, each time a new belief point is added, we select a subset of belief points that satisfy the approximate linear dependence (ALD) criterion to train the GPR. This selected subset is usually significantly smaller than the entire belief set, which further reduces computational demands.

Finally, it must be noted that the proposed GP-UCB is only approximately correct, although with a very high probability. Hence, we provide theoretical results on the consistency of the proposed upper bound approximations as the belief set expands.

B Description of the tiger problem

To demonstrate the concept of reachable belief points in a POMDP, we use the tiger problem, originally introduced by Kaelbling et al. (1998). In this problem, an agent faces two doors, with a tiger behind one and treasure behind the other, i.e., the unobservable states are “Tiger Left” and “Tiger Right.” At every stage, the agent can open one of the doors (i.e., actions “open left” or “open right”) or listen (action “listen”) for a tiger growl to gather (noisy) information. In particular, there is an 85% chance of correctly estimating the location of the tiger after listening. The agent’s actions are associated with specific rewards and costs: opening the correct door yields a reward of +10, while opening the door with the tiger results in a penalty of −100, and listening incurs a cost of −1. After the agent opens a door, the state resets: the location of the tiger is randomized, placing it behind one of the doors with equal probability. Observations are then made in the new state, where both left and right actions lead to either observation (Growl from Left or Growl from Right) with a probability of 0.5, regardless of the actual location of the tiger. In the finite-horizon version, at the terminating stage, only the immediate reward or penalty from the

final action is considered, without further actions or resets.

C Further details on sawtooth projection method

Hauskrecht (2000) proposed an interpolation approach for efficient approximation of the convex hull, referred to as *sawtooth projection*, to reduce computational complexity. Given an arbitrary belief $\mathbf{b} \in \mathcal{B}$ and corner beliefs $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{|\mathcal{S}|}\}$, the sawtooth projection provides the upper bound interpolation for any new belief point \mathbf{b}' . The corner beliefs are defined as $\mathbf{w}_s = (b(1), b(2), \dots, b(|\mathcal{S}|))$, where $b(s) = 1$ and $b(s') = 0$ for all $s' \neq s$. The best upper bound approximation for a belief \mathbf{b} is the one that minimizes $\lambda(\mathbf{b})f(\mathbf{b})$, where

$$\lambda(\mathbf{b}) = \min\{b'(s)/b(s) | s \in \mathcal{S}\}, \text{ and} \quad (\text{C.1})$$

$$f(\mathbf{b}) = \bar{V}(\mathbf{b}) - \sum_{s \in \mathcal{S}} b(s)\bar{V}(\mathbf{w}_s), \quad (\text{C.2})$$

Intuitively, this method interprets the value function as a set of downward-pointing pyramids, with their bases at the corners of the belief space and their tips at known belief point upper bounds.

Following the sawtooth projection for an arbitrary belief point \mathbf{b}' at time stage t , the upper bound for belief point, \mathbf{b} at stage t is obtained as

$$\bar{V}_t(\mathbf{b}) = \max_{a \in \mathcal{A}} \left[\sum_{s \in \mathcal{S}} b(s)R(s, a) + \sum_{o \in \mathcal{O}} P(o|\mathbf{b}, a)\bar{v}_{t+1}(\mathbf{b}') \right], \quad (\text{C.3})$$

where $\bar{v}_{t+1}(\mathbf{b}')$ is the upper bound projection obtained using sawtooth projection for the belief state \mathbf{b}' resulting from taking action a and observing o . Details of the upper bound projection using the sawtooth projection is presented in Algorithm C.1. Intuitively, the upper bound update is the maximum expected value of the sum of immediate reward and the upper bound of the future reward.

Algorithm C.1 Upperbound approximation with sawtooth projection method

- 1: **Input:** POMDP model, belief point \mathbf{b}' , belief-bound pairs $(\mathbf{b}, \bar{V}(\mathbf{b}))$ for each belief \mathbf{b} in the belief set \mathcal{B}
 - 2: **Output:** upper bound corresponding to belief point \mathbf{b}'
 - 3: **for** $\mathbf{b} \in \mathcal{B} \setminus \{\exists_f | f \in \mathcal{S}\}$ **do**
 - 4: $f(\mathbf{b}) \leftarrow \bar{V}(\mathbf{b}) - \sum_{s \in \mathcal{S}} b(s)\bar{V}(\mathbf{w}_s)$
 - 5: $\lambda(\mathbf{b}) \leftarrow \min_s \frac{b'(s)}{b(s)}$
 - 6: **end for**
 - 7: $\mathbf{b}^* \leftarrow \arg \min_{\mathbf{b}} \lambda(\mathbf{b})f(\mathbf{b})$
 - 8: **return** $\lambda(\mathbf{b}^*)f(\mathbf{b}^*) + \sum_{s \in \mathcal{S}} b'(s)\bar{V}(\mathbf{w}_s)$
-

Figure C.1 illustrates the sawtooth projection method for a two-state problem. The orange dots labeled 1 to 5 represent belief points with known upper bounds, with dots 1 and 5 being corner beliefs. By connecting belief points 2, 3, and 4 to the corner beliefs, three downward pointing

triangles are formed. For a new belief point \mathbf{b}' , the sawtooth method estimates its upper bound by projecting \mathbf{b}' onto the connecting lines and choosing the projection with the lowest value, which in this case is the line formed by dot 3 and corner belief 1. The estimated upper bound is represented by the red dot on the left of Figure C.1. The sawtooth projection for the entire belief space is shown as the red solid line in the right figure. For comparison, the optimal value function is plotted as the purple line in Figure C.1.

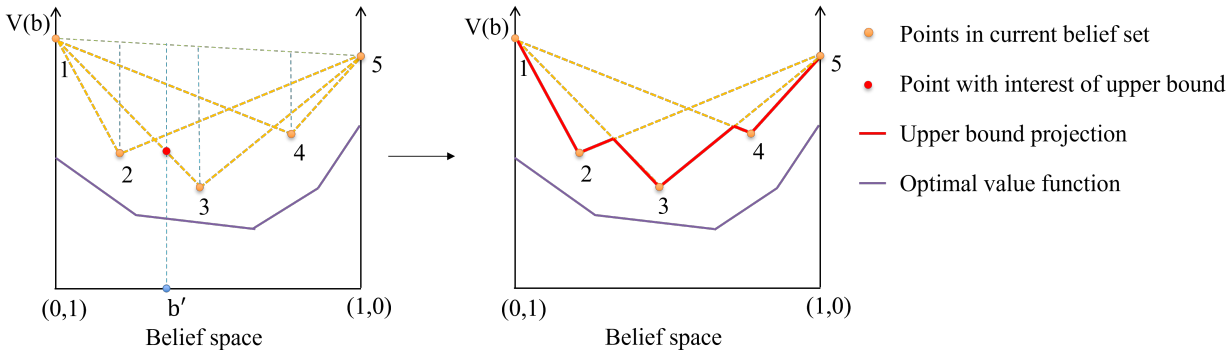


Figure C.1: The approximation of upper bound with sawtooth projection.

D Full computational results

Table D.1: Experiment results for the ChengD51 problem. The target gap is 0.001. Bolded numbers indicate the largest upper bounds and smallest gaps, values in parentheses are standard deviations.

	Max-gap GP-UCB	Max-gap sawtooth	Random GP-UCB	Random sawtooth	Fixed-grid GP-UCB	Fixed-grid sawtooth	
h=10	LB	65.246 ^(0.000)	65.246	65.246 ^(0.000)	65.246 ^(0.000)	65.246	
	Gap	0.004 ^(0.000)	0.005	0.231 ^(0.026)	0.263 ^(0.020)	16.897	18.896
	Time (s)	3008.2 ^(7.0)	3022.4	3013.8 ^(8.2)	3017.4 ^(6.9)	0.4	0.1
	Worst Gap	0.004	0.005	0.278	0.308	16.897	18.896
h=15	LB	97.990 ^(0.000)	97.990	97.990 ^(0.000)	97.990 ^(0.000)	97.990	97.990
	Gap	0.017 ^(0.000)	0.021	0.488 ^(0.014)	0.562 ^(0.049)	26.322	28.322
	Time (s)	3020.0 ^(17.6)	3010.9	3014.7 ^(10.1)	3018.2 ^(9.4)	0.7	0.2
	Worst Gap	0.017	0.021	0.502	0.650	26.322	28.322
h=20	LB	130.735 ^(0.000)	130.735	130.735 ^(0.000)	130.735 ^(0.000)	130.734	130.734
	Gap	0.034 ^(0.001)	0.045	0.758 ^(0.051)	0.850 ^(0.039)	35.748	37.747
	Time (s)	3019.9 ^(21.1)	3032.9	3047.9 ^(17.8)	3033.9 ^(8.5)	1.1	0.2
	Worst Gap	0.035	0.045	0.842	0.913	35.748	37.747
h=40	LB	261.713 ^(0.000)	261.713	261.713 ^(0.000)	261.713 ^(0.000)	261.711	261.711
	Gap	0.129 ^(0.002)	0.170	2.087 ^(0.058)	2.287 ^(0.092)	73.452	75.449
	Time (s)	3022.0 ^(9.3)	3028.8	3069.6 ^(26.9)	3019.3 ^(22.2)	2.3	0.5
	Worst Gap	0.131	0.170	2.170	2.447	73.452	75.449

Table D.2: Experiment results for the network problem. The target gap is 0.01. Bolded numbers indicate the largest upper bounds and smallest gaps, values in parentheses are standard deviations.

		Max-gap GP-UCB	Max-gap sawtooth	Random GP-UCB	Random sawtooth	Fixed-grid GP-UCB	Fixed-grid sawtooth
h=10	LB	151.180 ^(0.000)	151.180	151.180 ^(0.000)	151.180 ^(0.000)	150.456	150.456
	Gap	0.010 ^(0.000)	0.010	1.348 ^(0.232)	1.846 ^(0.357)	104.631	141.734
	Time (s)	58.8 ^(7.2)	52.9	3014.2 ^(11.9)	3012.2 ^(7.0)	0.7	0.3
	Worst Gap	0.010	0.010	1.983	2.365	104.631	141.734
h=15	LB	224.616 ^(0.000)	224.616	224.608 ^(0.009)	224.576 ^(0.079)	221.643	221.643
	Gap	0.025 ^(0.002)	0.039	5.888 ^(0.542)	8.076 ^(0.830)	176.103	229.092
	Time (s)	3011.7 ^(13.5)	3011.3	3018.5 ^(15.1)	3012.9 ^(10.4)	1.2	0.5
	Worst Gap	0.031	0.039	6.535	9.748	176.103	229.092
h=20	LB	298.149 ^(0.000)	298.149	298.101 ^(0.032)	298.053 ^(0.084)	293.685	293.685
	Gap	0.164 ^(0.024)	0.346	12.579 ^(1.042)	15.545 ^(0.995)	248.946	315.601
	Time (s)	3033.7 ^(20.6)	3029.9	3029.0 ^(21.6)	3011.9 ^(11.1)	1.7	0.7
	Worst Gap	0.216	0.346	13.854	17.795	248.946	315.601
h=40	LB	592.274 ^(0.003)	592.274	591.878 ^(0.150)	591.719 ^(0.229)	582.704	582.704
	Gap	1.453 ^(0.018)	2.419	42.844 ^(1.844)	50.222 ^(1.503)	543.012	660.788
	Time (s)	3030.6 ^(36.2)	3020.2	3015.0 ^(15.7)	3020.2 ^(13.1)	3.2	1.5
	Worst Gap	1.488	2.419	45.825	53.214	543.012	660.788

Table D.3: Experiment results for the query problem. The target gap is 0.001. Bolded numbers indicate the largest upper bounds and smallest gaps, values in parentheses are standard deviations. The fixed-grid method is unable to solve the problem within the limited 3000 seconds.

		Max-gap GP-UCB	Max-gap sawtooth	Random GP-UCB	Random sawtooth	Fixed-grid GP-UCB	Fixed-grid sawtooth
h=10	LB	30.021 ^(0.000)	30.021	30.021 ^(0.000)	30.021 ^(0.000)	NA	NA
	Gap	0.001 ^(0.000)	0.001	0.001 ^(0.000)	0.001 ^(0.000)	NA	NA
	Time (s)	5.7 ^(0.6)	80.7	7.3 ^(2.5)	3030.8 ^(7.9)	NA	NA
	Worst Gap	0.001	0.001	0.001	0.001	NA	NA
h=15	LB	45.047 ^(0.000)	45.047	45.047 ^(0.000)	45.047 ^(0.000)	NA	NA
	Gap	0.001 ^(0.000)	0.002	0.001 ^(0.000)	0.004 ^(0.000)	NA	NA
	Time (s)	8.6 ^(0.6)	3000.8	30.6 ^(8.4)	3036.2 ^(13.9)	NA	NA
	Worst Gap	0.001	0.002	0.001	0.004	NA	NA
h=20	LB	60.083 ^(0.000)	60.083	60.083 ^(0.000)	60.083 ^(0.000)	NA	NA
	Gap	0.001 ^(0.000)	0.004	0.001 ^(0.000)	0.007 ^(0.000)	NA	NA
	Time (s)	3007.6 ^(4.8)	3025.6	206.2 ^(72.9)	3039.5 ^(21.0)	NA	NA
	Worst Gap	0.001	0.004	0.001	0.007	NA	NA
h=40	LB	120.327 ^(0.000)	120.327	120.327 ^(0.000)	120.327 ^(0.000)	NA	NA
	Gap	0.017 ^(0.001)	0.026	0.010 ^(0.001)	0.037 ^(0.000)	NA	NA
	Time (s)	3004.4 ^(5.6)	3049.5	3018.8 ^(32.9)	3027.4 ^(30.6)	NA	NA
	Worst Gap	0.018	0.026	0.010	0.037	NA	NA

Table D.4: Experiment results for the hallway problem. The target gap is 0.00001. Bolded numbers indicate the largest upper bounds and smallest gaps, values in parentheses are standard deviations. The fixed-grid method is unable to solve the problem within the limited 3000 seconds.

		Max-gap GP-UCB	Max-gap sawtooth	Random GP-UCB	Random sawtooth	Fixed-grid GP-UCB	Fixed-grid sawtooth
h=10	LB	0.311 ^(0.004)	0.311	0.327 ^(0.002)	0.326 ^(0.003)	NA	NA
	Gap	0.162 ^(0.000)	0.168	0.188 ^(0.005)	0.180 ^(0.009)	NA	NA
	Time (s)	3067.0 ^(27.9)	3034.5	3054.3 ^(36.8)	3065.6 ^(47.5)	NA	NA
	Worst Gap	0.162	0.168	0.196	0.196	NA	NA
h=15	LB	0.604 ^(0.001)	0.587	0.591 ^(0.012)	0.587 ^(0.011)	NA	NA
	Gap	0.334 ^(0.001)	0.360	0.384 ^(0.013)	0.387 ^(0.010)	NA	NA
	Time (s)	2989.8 ^(24.5)	3082.8	3086.2 ^(57.0)	3062.6 ^(33.9)	NA	NA
	Worst Gap	0.336	0.360	0.396	0.391	NA	NA
h=20	LB	0.865 ^(0.001)	0.858	0.811 ^(0.024)	0.796 ^(0.023)	NA	NA
	Gap	0.507 ^(0.001)	0.524	0.604 ^(0.030)	0.615 ^(0.026)	NA	NA
	Time (s)	3022.0 ^(25.1)	2883.4	3010.1 ^(37.3)	3063.4 ^(41.6)	NA	NA
	Worst Gap	0.508	0.524	0.646	0.635	NA	NA
h=40	LB	1.930 ^(0.000)	1.912	1.612 ^(0.067)	1.295 ^(0.061)	NA	NA
	Gap	1.220 ^(0.000)	1.234	1.451 ^(0.070)	1.510 ^(0.061)	NA	NA
	Time (s)	2871.8 ^(18.9)	3030.7	3092.7 ^(66.6)	3069.2 ^(21.0)	NA	NA
	Worst Gap	1.220	1.234	1.571	1.638	NA	NA

Table D.5: Experiment results for the aloha30 problem. The target gap is 0.01. Bolded numbers indicate the largest upper bounds and smallest gaps, values in parentheses are standard deviations. The fixed-grid method is unable to solve the problem within the limited 3000 seconds.

		Max-gap GP-UCB	Max-gap sawtooth	Random GP-UCB	Random sawtooth	Fixed-grid GP-UCB	Fixed-grid sawtooth
h=10	LB	122.527 ^(0.000)	122.501	122.545 ^(0.006)	122.517 ^(0.016)	NA	NA
	Gap	0.735 ^(0.000)	0.736	0.877 ^(0.072)	0.963 ^(0.121)	NA	NA
	Time (s)	3042.3 ^(6.2)	2962.7	3035.4 ^(21.3)	3004.7 ^(22.7)	NA	NA
	Worst Gap	0.735	0.736	0.979	1.177	NA	NA
h=15	LB	167.681 ^(0.000)	167.591	167.784 ^(0.017)	167.625 ^(0.137)	NA	NA
	Gap	1.836 ^(0.000)	2.004	2.026 ^(0.132)	2.460 ^(0.229)	NA	NA
	Time (s)	3050.9 ^(21.2)	2912.5	3034.1 ^(29.2)	3101.9 ^(77.6)	NA	NA
	Worst Gap	1.836	2.004	2.221	3.017	NA	NA
h=20	LB	203.827 ^(0.004)	203.657	204.090 ^(0.045)	203.765 ^(0.203)	NA	NA
	Gap	3.248 ^(0.007)	4.003	3.661 ^(0.095)	4.417 ^(0.317)	NA	NA
	Time (s)	3061.5 ^(12.4)	3070.2	3073.9 ^(32.9)	3134.4 ^(33.9)	NA	NA
	Worst Gap	3.262	4.003	3.804	5.134	NA	NA
h=40	LB	282.129 ^(0.000)	282.039	283.338 ^(0.249)	282.329 ^(0.577)	NA	NA
	Gap	10.428 ^(0.003)	12.507	11.164 ^(0.632)	12.480 ^(0.932)	NA	NA
	Time (s)	2956.0 ^(27.2)	3086.8	3084.9 ^(45.8)	3163.7 ^(39.4)	NA	NA
	Worst Gap	10.435	12.507	12.113	14.570	NA	NA

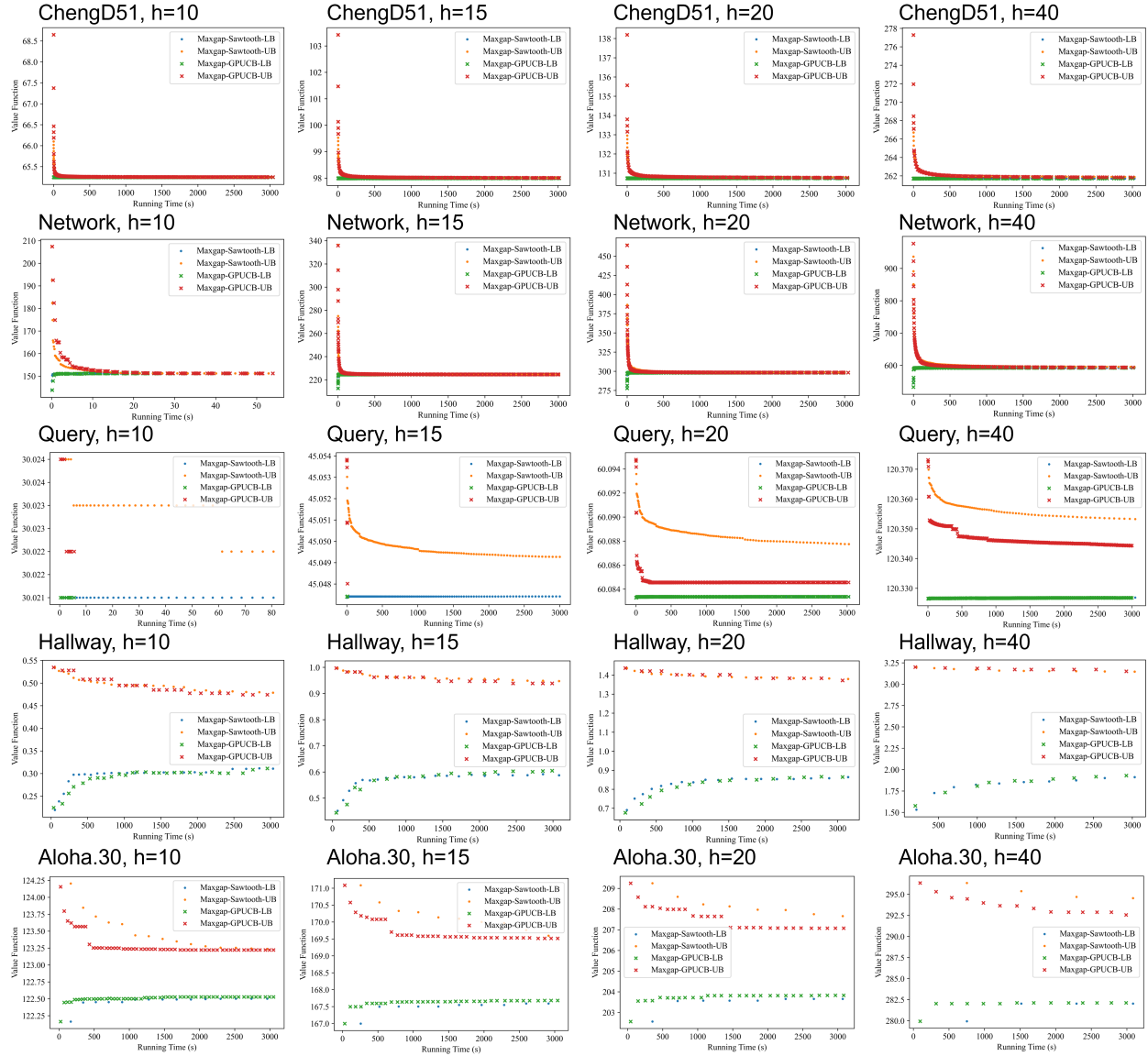


Figure D.1: Upper and lower bounds of the value function for the five examples. Dots represent upper and lower bounds from FiVI, while crosses indicate those from AUCB, color-coded as shown in the legend.