

Simplifying Complex Observation Models in Continuous POMDP Planning with Probabilistic Guarantees and Practice

Idan Lev-Yehudi Moran Barenboim Vadim Indelman



ANPL
Autonomous Navigation and
Perception Lab



TECHNION
Israel Institute
of Technology

Introduction

- ▶ Planning under uncertainty can be formalized as a Partially Observable Markov Decision Process (POMDP)

¹Wang et al., “DualSMC: Tunneling Differentiable Filtering and Planning under Continuous POMDPs”; Deglurkar et al., “Compositional Learning-based Planning for Vision POMDPs”.

Introduction

- ▶ Planning under uncertainty can be formalized as a Partially Observable Markov Decision Process (POMDP)
- ▶ Optimally solving POMDPs is computationally expensive and feasible only for small tasks

¹Wang et al., “DualSMC: Tunneling Differentiable Filtering and Planning under Continuous POMDPs”; Deglurkar et al., “Compositional Learning-based Planning for Vision POMDPs”.

Introduction

- ▶ Planning under uncertainty can be formalized as a Partially Observable Markov Decision Process (POMDP)
- ▶ Optimally solving POMDPs is computationally expensive and feasible only for small tasks
- ▶ Visual observations are complex to model in planning¹

¹Wang et al., “DualSMC: Tunneling Differentiable Filtering and Planning under Continuous POMDPs”; Deglurkar et al., “Compositional Learning-based Planning for Vision POMDPs”.

Introduction

- ▶ Planning under uncertainty can be formalized as a Partially Observable Markov Decision Process (POMDP)
- ▶ Optimally solving POMDPs is computationally expensive and feasible only for small tasks
- ▶ Visual observations are complex to model in planning¹
- ▶ Learned observation models are impractical for solving the POMDP in real-time

¹Wang et al., “DualSMC: Tunneling Differentiable Filtering and Planning under Continuous POMDPs”; Deglurkar et al., “Compositional Learning-based Planning for Vision POMDPs”.

Contribution

- ▶ We explore planning with a simpler observation model while attaining formal guarantees of the solution quality

Contribution

- ▶ We explore planning with a simpler observation model while attaining formal guarantees of the solution quality
- ▶ Potential of substantial computational improvement for complex models

Contribution

- ▶ We explore planning with a simpler observation model while attaining formal guarantees of the solution quality
- ▶ Potential of substantial computational improvement for complex models
- ▶ Our main contributions:
 - ▶ Bound the theoretical loss with observation model discrepancy
 - ▶ Probabilistic bound for the empirical simplified performance
 - ▶ Practical computation of the bounds in SOTA planners

Continuous POMDP Solvers

- ▶ POMCPOW is a SOTA continuous POMDP solver ²

Algorithm 2 POMCPOW

```
1: procedure SIMULATE( $s, h, d$ )
2:   if  $d = 0$  then
3:     return 0
4:    $a \leftarrow \text{ACTIONPROGWIDEN}(h)$ 
5:    $s', o, r \leftarrow G(s, a)$ 
6:   if  $|C(ha)| \leq k_o N(ha)^{\alpha_o}$  then
7:      $M(hao) \leftarrow M(hao) + 1$ 
8:   else
9:      $o \leftarrow \text{select } o \in C(ha) \text{ w.p. } \frac{M(hao)}{\sum_o M(hao)}$ 
10:   $\text{append } s' \text{ to } B(hao)$ 
11:   $\text{append } \mathcal{Z}(o \mid s, a, s') \text{ to } W(hao)$ 
12:  if  $o \notin C(ha)$  then ▷ new node
13:     $C(ha) \leftarrow C(ha) \cup \{o\}$ 
14:     $total \leftarrow r + \gamma \text{ROLLOUT}(s', hao, d - 1)$ 
15:  else
16:     $s' \leftarrow \text{select } B(hao)[i] \text{ w.p. } \frac{W(hao)[i]}{\sum_{j=1}^m W(hao)[j]}$ 
17:     $r \leftarrow R(s, a, s')$ 
18:     $total \leftarrow r + \gamma \text{SIMULATE}(s', hao, d - 1)$ 
19:   $N(h) \leftarrow N(h) + 1$ 
20:   $N(ha) \leftarrow N(ha) + 1$ 
21:   $Q(ha) \leftarrow Q(ha) + \frac{total - Q(ha)}{N(ha)}$ 
22:  return  $total$ 
```

²Sunberg and Kochenderfer, "Online algorithms for POMDPs with continuous state, action, and observation spaces"

Problem Formulation

- ▶ A POMDP is the tuple $\langle \mathcal{X}, \mathcal{A}, \mathcal{Z}, p_T, p_Z, r, \gamma, L, b_0 \rangle$
 - ▶ $\mathcal{X}, \mathcal{A}, \mathcal{Z}$ are state, action and observation spaces
 - ▶ p_T, p_Z are probabilistic transition and observation models
 - ▶ $r_t: \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ is a bounded reward function at time t
 - ▶ γ is the reward discount for future time steps
 - ▶ L is the time limit (horizon)
 - ▶ b_0 is the starting distribution (belief) of states







Problem Formulation

- ▶ A POMDP is the tuple $\langle \mathcal{X}, \mathcal{A}, \mathcal{Z}, p_T, p_Z, r, \gamma, L, b_0 \rangle$
 - ▶ $\mathcal{X}, \mathcal{A}, \mathcal{Z}$ are state, action and observation spaces
 - ▶ p_T, p_Z are probabilistic transition and observation models
 - ▶ $r_t: \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ is a bounded reward function at time t
 - ▶ γ is the reward discount for future time steps
 - ▶ L is the time limit (horizon)
 - ▶ b_0 is the starting distribution (belief) of states
- ▶ Action-value function:

$$Q_{\mathbf{P}}^{p_Z}(b_t, a) \triangleq r_t(b_t, a) + \mathbb{E}_{z_{t+1:L} \sim p_Z} [\sum_{i=t+1}^L \gamma^{i-t} r_i(b_i, \pi_i)]$$







Simplifying the Observation Model

- ▶ We replace p_Z with a cheaper model q_Z
- ▶ Simplified Action-value function: $Q_P^{q_Z}$

	Theoretical Values	PB-MDP	Planner
Original Model			
Simplified Model			

Simplifying the Observation Model

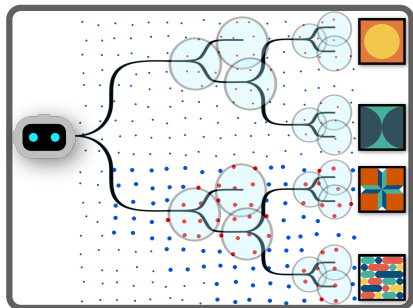
- ▶ We replace p_Z with a cheaper model q_Z
- ▶ Simplified Action-value function: $Q_P^{q_Z}$

	Theoretical Values	PB-MDP	Planner
Original Model			
Simplified Model			

- ▶ Can we bound $|Q_P^{p_Z} - \hat{Q}_{M_P}^{q_Z}|$?

Approach to Bounds

- ▶ Pre-sample states at which we compute "observation model discrepancy"
- ▶ During online, we weight these states according to their likelihood
- ▶ We prove convergence guarantees for our estimated bounds



State-Dependent Observation TV-Distance

- ▶ Obs. TV-Distance: $\Delta_Z(x) \triangleq \int_{\mathcal{Z}} |p_Z(z | x) - q_Z(z | x)| dz$

State-Dependent Observation TV-Distance

- ▶ Obs. TV-Distance: $\Delta_Z(x) \triangleq \int_{\mathcal{Z}} |p_Z(z | x) - q_Z(z | x)| dz$
- ▶ $m_i(x_i, a) \triangleq V_{i+1}^{\max} \cdot \mathbb{E}_{x_{i+1} \sim p_T(\cdot | x_i, a)} [\Delta_Z(x_{i+1})]$

State-Dependent Observation TV-Distance

- ▶ Obs. TV-Distance: $\Delta_Z(x) \triangleq \int_{\mathcal{Z}} |p_Z(z | x) - q_Z(z | x)| dz$
- ▶ $m_i(x_i, a) \triangleq V_{i+1}^{\max} \cdot \mathbb{E}_{x_{i+1} \sim p_T(\cdot | x_i, a)} [\Delta_Z(x_{i+1})]$
- ▶ $m_i(b_i, a) \triangleq \mathbb{E}_{x_i \sim b_i} [m_i(x_i, a)]$

State-Dependent Observation TV-Distance

- ▶ Obs. TV-Distance: $\Delta_Z(x) \triangleq \int_{\mathcal{Z}} |p_Z(z | x) - q_Z(z | x)| dz$
- ▶ $m_i(x_i, a) \triangleq V_{i+1}^{\max} \cdot \mathbb{E}_{x_{i+1} \sim p_T(\cdot | x_i, a)} [\Delta_Z(x_{i+1})]$
- ▶ $m_i(b_i, a) \triangleq \mathbb{E}_{x_i \sim b_i} [m_i(x_i, a)]$
- ▶ It is natural to define cumulative bound function

$$\Phi_{\mathbf{P}}(b_t, a) \triangleq m_t(b_t, a) + \mathbb{E}_{z_{t+1:L-1} \sim q_Z} \left[\sum_{i=t+1}^{L-1} m_i(b_i, \pi_i) \right]$$

$$Q_{\mathbf{P}}^{p_Z}(b_t, a) \triangleq r_t(b_t, a) + \mathbb{E}_{z_{t+1:L} \sim p_Z} \left[\sum_{i=t+1}^L \gamma^{i-t} r_i(b_i, \pi_i) \right]$$

TV-Distance Loss Bounds

Theorem 2

For every belief b_t , action a , policy π , observation models p_Z and q_Z , the following bound holds deterministically:

$$|Q_{\mathbf{P}}^{p_Z}(b_t, a) - Q_{\mathbf{P}}^{q_Z}(b_t, a)| \leq \Phi_{\mathbf{P}}(b_t, a)$$

TV-Distance Loss Bounds

Theorem 2

For every belief b_t , action a , policy π , observation models p_Z and q_Z , the following bound holds deterministically:

$$|Q_{\mathbf{P}}^{p_Z}(b_t, a) - Q_{\mathbf{P}}^{q_Z}(b_t, a)| \leq \Phi_{\mathbf{P}}(b_t, a)$$

Theorem 3 (Informal)

For every bounded state-action function (r_i/m_i) , its finite-sample cumulative function $(Q_{\mathbf{M}_{\mathbf{P}}}^{q_Z}/\Phi_{\mathbf{M}_{\mathbf{P}}})$ has probabilistic concentration bounds from its theoretical counterpart $(Q_{\mathbf{P}}^{q_Z}/\Phi_{\mathbf{P}})$ under certain regularity conditions of the POMDP

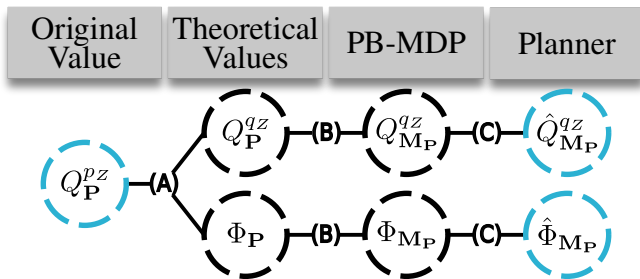
Empirical Concentration Inequalities

Corollary 3

For arbitrary $\varepsilon, \delta > 0$ there exists a number of particles for which

$$|Q_{\mathbf{P}}^{pZ}(b_t, a) - \hat{Q}_{\mathbf{M}_{\mathbf{P}}}^{qZ}(\bar{b}_t, a)| \leq \hat{\Phi}_{\mathbf{M}_{\mathbf{P}}}(\bar{b}_t, a) + \varepsilon$$

with probability of at least $1 - \delta$ for any guaranteed planner



- ▶ (A) is given by Theorem 2, (B) is given by Theorem 3, (C) is given by any planner with performance guarantees

Practical Computation of Bounds

- ▶ Computation of m_i is impractical
 - ▶ Importance Sampling
 - ▶ Separate calculations to offline/online

Online

$$\tilde{m}_i(x_i, a) \triangleq V_{i+1}^{\max} \frac{1}{N_{\Delta}} \sum_{n=1}^{N_{\Delta}} \frac{p_T(x_n^{\Delta} | x_i, a)}{Q_0(x_n^{\Delta})} \Delta_Z(x_n^{\Delta})$$
$$\{x_n^{\Delta}\}_{n=1}^{N_{\Delta}} \sim Q_0(x)$$

Offline

Practical Computation of Bounds

- ▶ Computation of m_i is impractical
 - ▶ Importance Sampling
 - ▶ Separate calculations to offline/online

The diagram illustrates the flow of information between Online and Offline components. A grey box labeled "Online" has two arrows pointing to the terms $\frac{1}{N_\Delta} \sum_{i=1}^{N_\Delta} \frac{p_T(x_n^\Delta | x_i, a)}{Q_0(x_n^\Delta)}$ and $\Delta_Z(x_n^\Delta)$ in the equation. A grey box labeled "Offline" has two arrows pointing to the terms $\frac{p_T(x_n^\Delta | x_i, a)}{Q_0(x_n^\Delta)}$ and $\{x_n^\Delta\}_{n=1}^{N_\Delta} \sim Q_0(x)$ in the equation.

$$\tilde{m}_i(x_i, a) \triangleq V_{i+1}^{\max} \frac{1}{N_\Delta} \sum_{i=1}^{N_\Delta} \frac{p_T(x_n^\Delta | x_i, a)}{Q_0(x_n^\Delta)} \Delta_Z(x_n^\Delta)$$
$$\{x_n^\Delta\}_{n=1}^{N_\Delta} \sim Q_0(x)$$

- ▶ Optimizations:
 - ▶ Considering state-samples based on a KD-Tree and a truncation distance
 - ▶ Computing a Monte Carlo estimate of \tilde{m}_i .

Practical Computation of Bounds

- ▶ Computation of m_i is impractical
 - ▶ Importance Sampling
 - ▶ Separate calculations to offline/online

Online

$$\tilde{m}_i(x_i, a) \triangleq V_{i+1}^{\max} \frac{1}{N_{\Delta}} \sum_{n=1}^{N_{\Delta}} \frac{p_T(x_n^{\Delta} | x_i, a)}{Q_0(x_n^{\Delta})} \Delta_Z(x_n^{\Delta})$$

Offline

$$\{x_n^{\Delta}\}_{n=1}^{N_{\Delta}} \sim Q_0(x)$$

- ▶ Optimizations:
 - ▶ Considering state-samples based on a KD-Tree and a truncation distance
 - ▶ Computing a Monte Carlo estimate of \tilde{m}_i .
- ▶ In the paper we discuss embedding \tilde{m}_i into POMDP solvers

Results in Simulation

- ▶ We show in our simulative setup that even with bounds calculation we achieve a significant speedup



Deglurkar, Sampada et al. “Compositional Learning-based Planning for Vision POMDPs”. In: *Learning for Dynamics and Control Conference*. PMLR. 2023, pp. 469–482.



Sunberg, Zachary and Mykel Kochenderfer. “Online algorithms for POMDPs with continuous state, action, and observation spaces”. In: *Proceedings of the International Conference on Automated Planning and Scheduling*. Vol. 28. 1. 2018.



Wang, Yunbo et al. “DualSMC: Tunneling Differentiable Filtering and Planning under Continuous POMDPs”. In: *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*. Ed. by Christian Bessiere. Main track. International Joint Conferences on Artificial Intelligence Organization, July 2020, pp. 4190–4198. DOI: 10.24963/ijcai.2020/579. URL: <https://doi.org/10.24963/ijcai.2020/579>.