

Planning in Partially Observable Markov Decision Process

Shvechikov Pavel

National Research University Higher School of Economics,
Yandex School of Data Analysis

May 11, 2017

Overview

- 1 POMDP framework
 - POMDP model
- 2 Approximate planning
 - Exact planning – value iteration
 - Online vs offline planning
 - Online planning: POMCP

POMDP place in a model world

Markov Models		Do we have control over the state transitions?	
		NO	YES
Are the states completely observable?	YES	Markov Chain	MDP Markov Decision Process
	NO	HMM Hidden Markov Model	POMDP Partially Observable Markov Decision Process

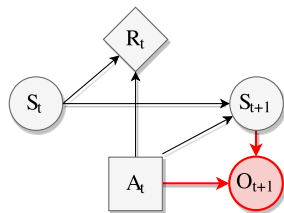
POMDP model relations

POMDP model

Definition

Partially Observed Markov Decision Process is a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \Omega, \mathcal{O} \rangle$

- 1 $\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}$ are the same as in MDP
- 2 Ω – finite set of observations
- 3 $\mathcal{O} : \mathcal{S} \times \mathcal{A} \mapsto \Delta(\Omega)$ – observation function, which gives, for each state and action, a probability distribution over Ω , i.e. $p(o | s_{t+1}, a_t) \quad \forall o \in \Omega$







Reasoning about state uncertainty



Belief state

Distribution over state space, i.e. $\sum_{s \in \mathcal{S}} b(s) = 1, \quad 0 \leq b(s) \leq 1$

$$\mathcal{A} = \{left, right\}, \quad p(\bar{A} | do(A)) = 0.1$$



$b(s_1)$	$b(s_2)$	$b(s_3)$	$b(s_4)$
0.333	0.333		0.333
0.1	0.450		0.450
0.1	0.164		0.736



Belief updating (Bayes filter)

Good news: belief updating is rather straightforward (Bayes Rule)

$$\begin{aligned} b'(s') &= p(s' | o', a, b) = \frac{p(o' | s', a) \cdot p(s' | a, b)}{\sum_o p(o' | s', a) \cdot p(s' | a, b)} \\ &\propto p(o' | s', a) \cdot p(s' | a, b) \\ &\propto p(o' | s', a) \sum_s p(s' | a, b, s) \cdot p(s | a, b) \\ &\propto p(o' | s', a) \sum_s p(s' | a, s) \cdot b(s) \end{aligned}$$

Bad news: belief updating can be computed exactly only for

- ① discrete low-dimensional state-spaces
- ② linear-Gaussian dynamics (leading to Kalman filter), i.e.
 - $s' \sim \mathcal{N}(s' | T_s s + T_a a, \Sigma_s), \quad o' \sim \mathcal{N}(o' | O_s s', \Sigma_o)$
 - $R(s, a) = s^\top R_s s + a^\top R_a a$

Discrete Bayes filter: **particle filter**

Basic idea: represent $b(\cdot)$ as set of particles $b^i \in \mathcal{S}$, $i = 1, \dots, K$

Particle filter with rejection

- ➊ **Givens:** $b = \{b_t^i \mid i = 1, \dots, K\}$, a , o
- ➋ **Set:** $b' = \{\emptyset\}$
- ➌ **Repeat** K times
 - ➊ Sample random state s from b
 - ➋ Repeat $(s', o') \sim G(s, a)$ until $o' = o$
 - ➌ Add s' to b'
- ➍ **Return:** b'

Depletion of a particle set is solved with reinvigoration:

- ➊ addition of random particles
- ➋ perturbation of some particles to nearby states

Outline

- 1 POMDP framework
 - POMDP model
- 2 Approximate planning
 - Exact planning – value iteration
 - Online vs offline planning
 - Online planning: POMCP

POMDP planning

We **know**:

- System dynamics, i.e. $p(s_{t+1} | s_t, a)$
- Observation function, i.e. $p(o | s_{t+1}, a)$
- Reward function, i.e. $\mathcal{R}(s_t, a)$

We **want**:

- $\max_{\pi} V_{\pi}(b_0) = \sum_s b_0(s) V_{\pi}(s) = \sum_s b_0(s) \cdot \mathbb{E} [\sum_t \gamma^t R_t | \pi]$

POMDP planning

We **know**:

- System dynamics, i.e. $p(s_{t+1} | s_t, a)$
- Observation function, i.e. $p(o | s_{t+1}, a)$
- Reward function, i.e. $\mathcal{R}(s_t, a)$

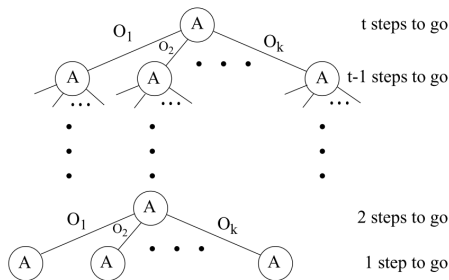
We **want**:

- $\max_{\pi} V_{\pi}(b_0) = \sum_s b_0(s) V_{\pi}(s) = \sum_s b_0(s) \cdot \mathbb{E} [\sum_t \gamma^t R_t | \pi]$

Number of policies to consider is

$$|\mathcal{A}| \sum_{t=0}^{T-1} |\mathcal{O}|^t = |\mathcal{A}| \frac{|\mathcal{O}|^T - 1}{|\mathcal{O}| - 1}$$

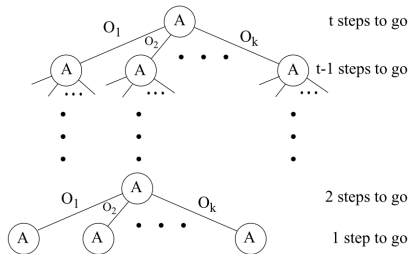
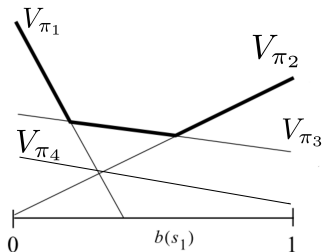
$|\mathcal{A}| = 2$, $|\mathcal{O}| = 2$, $T = 10$
yields 2^{1023} policies



Intro in exact planning in POMDP

$$V_{\pi}(s) = R(s, \pi(s)) + \gamma \sum_{s'} p(s' | s, \pi(s)) \sum_o p(o | s', \pi(s)) V_{o'(\pi)}(s')$$

$$V_{\pi}(\mathbf{b}) = \sum_s V_{\pi}(s) b(s) = [\mathbf{V}_{\pi}(s_1) \dots \mathbf{V}_{\pi}(s_n)] \cdot [\mathbf{b}(s_1) \dots \mathbf{b}(s_n)]^{\top} = \boldsymbol{\alpha}_{\pi} \cdot \mathbf{b}^{\top}$$



POMDP exact planning: value iteration

To sum up, we know

- 1 how to compute $V(b)$
- 2 that POMDP is an MDP over belief states
- 3 *value iteration* algorithm for MDP

POMDP value iteration update

$$V_{t+1}(b) = \max_a \left[\sum_s b(s) R(s, a) + \gamma \sum_{b'} p(b' | a, b) V_t(b') \right]$$

$$p(b' | a, b) = \sum_{o'} p(b' | a, b, o') \sum_{s'} p(o' | s', a) \sum_s p(s' | s, a) b(s)$$

POMDP exact planning: value iteration

To sum up, we know

- 1 how to compute $V(b)$
- 2 that POMDP is an MDP over **belief states**
- 3 *value iteration* algorithm for MDP

POMDP value iteration update

$$V_{t+1}(b) = \max_a \left[\sum_s b(s) R(s, a) + \gamma \sum_{b'} p(b' | a, b) V_t(b') \right]$$

$$p(b' | a, b) = \sum_{o'} p(b' | a, b, o') \sum_{s'} p(o' | s', a) \sum_s p(s' | s, a) b(s)$$

PSPACE-complete in the worst case

Online vs offline planning

Planning under partial observability consists of:

- *exact* planning
- offline *approximate* planning
- online *approximate* planning

Factor	Online	Offline
Plans for	current belief	all possible beliefs
Most computation	during execution	prior to execution
$ \mathcal{S} , \mathcal{A} $ size	billions and more	thousands
Execution speed	slow	fast
Recent SOTA	DESPOT, POMCP	PLEASE, SARSOP

Partially Observable Monte-Carlo Planning (POMCP)

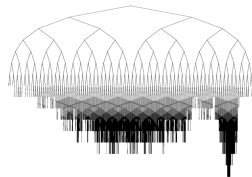
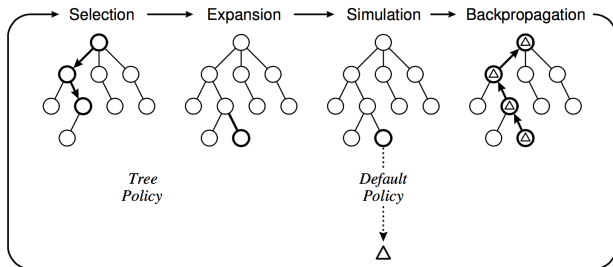
POMCP (Silver et al., 2010) ingredients:

- ① Monte-Carlo tree search (MCTS)
- ② Upper Confidence Bound policy (UCB1)
- ③ Particle filter

Key features:

- ① Requires **only generative model** $(s', o, r) \sim G(s, a)$
- ② **MC for both** belief updates and planning
- ③ Algorithm is **anytime**
- ④ Partially observed planning **at scale**

MCTS + UCB1 = UCT (Upper Confidence Tree)



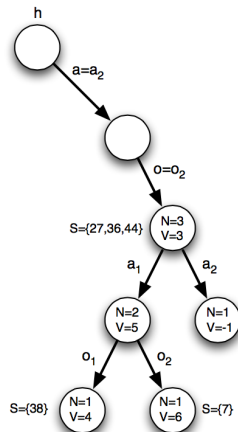
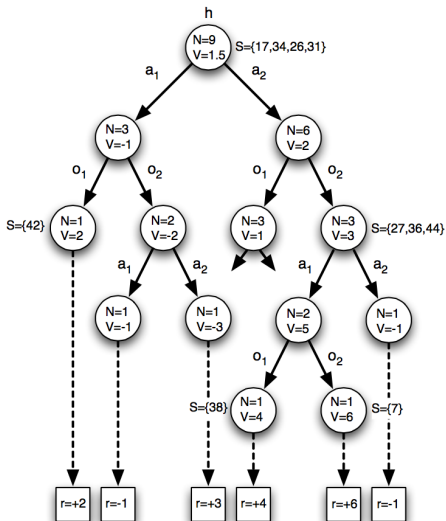
UCTB1 for multi-armed bandit policy maximizes $\bar{X}_j + \sqrt{\frac{\log n}{n_j}}$

UCTB1 applied to nodes in a tree gives UCT policy:

- ① If \exists **untried** a in s , expand s
- ② else greedy w.r.t. $Q^+(s, a) = Q(s, a) + c\sqrt{\frac{\log N(s)}{N(s, a)}}$

UCT + Particle filter = PO-UCT

Nodes in a tree are **histories**: i.e. $h = a_1 o_1 a_2 o_2 \dots a_n o_n$



POMCP

Algorithm 1 Partially Observable Monte-Carlo Planning

```

procedure SEARCH( $h$ )
  repeat
    if  $h = \text{empty}$  then
       $s \sim \mathcal{I}$ 
    else
       $s \sim B(h)$ 
    end if
    SIMULATE( $s, h, 0$ )
  until TIMEOUT()
  return  $\underset{b}{\operatorname{argmax}} V(hb)$ 
end procedure

procedure ROLLOUT( $s, h, \text{depth}$ )
  if  $\gamma^{\text{depth}} < \epsilon$  then
    return 0
  end if
   $a \sim \pi_{\text{rollout}}(h, \cdot)$ 
   $(s', o, r) \sim \mathcal{G}(s, a)$ 
  return  $r + \gamma \cdot \text{ROLLOUT}(s', hao, \text{depth}+1)$ 
end procedure

procedure SIMULATE( $s, h, \text{depth}$ )
  if  $\gamma^{\text{depth}} < \epsilon$  then
    return 0
  end if
  if  $h \notin T$  then
    for all  $a \in \mathcal{A}$  do
       $T(ha) \leftarrow (N_{\text{init}}(ha), V_{\text{init}}(ha), \emptyset)$ 
    end for
    return ROLLOUT( $s, h, \text{depth}$ )
  end if
   $a \leftarrow \underset{b}{\operatorname{argmax}} V(hb) + c \sqrt{\frac{\log N(h)}{N(hb)}}$ 
   $(s', o, r) \sim \mathcal{G}(s, a)$ 
   $R \leftarrow r + \gamma \cdot \text{SIMULATE}(s', hao, \text{depth}+1)$ 
   $B(h) \leftarrow B(h) \cup \{s\}$ 
   $N(h) \leftarrow N(h) + 1$ 
   $N(ha) \leftarrow N(ha) + 1$ 
   $V(ha) \leftarrow V(ha) + \frac{R - V(ha)}{N(ha)}$ 
  return  $R$ 
end procedure

```

Thank you!

References I



Silver, David et al. (2010). “Monte-Carlo planning in large POMDPs”. In: *Advances in neural information processing systems*, pp. 2164–2172.