

Last Time

Offline POMDP Algorithms

PBVI

SARSOP

$$\alpha_\pi[s]$$

Index - state
Numbers - Value of execution
Correspond - π at state s

$$V^\pi(b) = \alpha_\pi^T b$$

This Time

Formulation Approximations

Object Oriented Programming in Julia in HW5

Disappointing Facts:

Infinite Horizon POMDPs are Undecidable

Finite Horizon POMDPs are PSPACE complete.

Among the hardest problems that
can be solved using a polynomial amount
of space

~~Any~~.

POMDP algorithms (likely) have exponential complexity.

Numerical Approximation

Offline

Online

Solve the original problem approximately

Formulation Approximations

Solve an approximate
problem exactly
or approximately

Name	Description	Properties	Usefulness
Certainty - Equivalence	control as if the true state is mean of belief	Optimal for LQG	★★★★☆
QMDP	Full observability after 1 time step	QMDP pseudo α 's are upper bound for true V	★★★★☆
FIB	Takes 1 observation into account	tighter upper bound than QMDP	★★
Hindsight Opt ^{HOP}	Hindsight knowledge of state and outcome uncertainty	Looser Upper Bound than QMDP	★★★★
Last k observations "k Markov"	pretend that last k observations make up the state and solve that MDP	Great for Atari!	★★★★
Open Loop no observations	Choose sequence of actions	Good if Aleatory is low, and epistemic is hard to reduce	★★★
Most likely obs	Plan assuming $b' = \tau(b, a, \hat{o}(b))$	No observation branching	★★★

POMDP Objective

$$\pi^* = \operatorname{argmax}_{\pi: \mathcal{B} \rightarrow \mathcal{A}} E \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

Certainty Equivalence

$$\pi_{CE}(b) = \pi_{MDP}(E[s]_{s=b}) \quad \begin{array}{l} \text{or mode} \\ \text{or median} \end{array}$$

Optimal for one very special POMDP

LQG

$$\begin{aligned} s' &\sim N(A s + B a, V) \\ o &\sim N(C s' + D u, W) \end{aligned}$$

$$\begin{aligned} s_0 &\sim N(\mu_0, \Sigma_0) \\ R(s, a) &= s^T Q s + a^T R a \end{aligned}$$

Bayesian
Update

Kalman Filter
 $b_t = N(\mu_t, \Sigma_t)$

Kalman
Gain \rightarrow

$$\begin{aligned} \Sigma_{t+1} &= A (\Sigma_t - \Sigma_t C^T (C \Sigma_t C^T + W)^{-1} C \Sigma_t) A^T + V \\ L &= A \Sigma_t C^T (C \Sigma_t C^T + W)^{-1} \\ \mu_{t+1} &= A \mu_t + B a + L (o_t - \underline{C} \mu_t) \end{aligned}$$

Solution to LQR MDP

$$a_t = -K s_t$$

Solution to LQG POMDP

$$a_t = -K \mu_t$$

Certainty Equivalence Principle

Works pretty well for any problem where
belief is unimodal

QMDP

$$\alpha_a[s] = Q_{MDP}^*(s, a)$$

$$\pi_{QMDP}(b) = \operatorname{argmax}_{sub} E [Q_{MDP}(s, a)] \Rightarrow \operatorname{argmax}_a \alpha_a^T b$$

Breakout Rooms

$$b = U(\{1, 2\}) \quad \pi^*(b) = +10$$

What is $\pi_{QMDP}(b)$

Steps

figure out $Q_{MDP}^*(s, a)$

What would you do if you knew the state

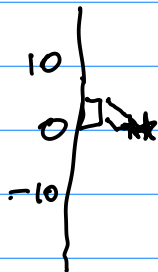
$$\operatorname{argmax}_{sub} E [Q_{MDP}(s, a)]$$

$$Q_{MDP}^*(1, -1) = 100 - 1$$

$$Q_{MDP}^*(2, -1) = 100 - 2$$

$$Q_{MDP}^*(1, a \neq -1) = \text{worse than } -1$$

$$Q_{MDP}^*(2, a \neq -1) = \text{worse than } -1$$



QMDP is bad at costly information gathering

FIB

$$\pi_{FIB}(b) = \operatorname{argmax}_a \alpha_a^T b$$

loop

$$\alpha_a^{(k+1)}[s] = R(s, a) + \gamma \sum_0 \max_{a'} \sum_{s'} \underbrace{Z(0|a, s')}_{\text{}} \underbrace{T(s'|s, a)}_{\text{}} \alpha_{a'}^{(k)}[s']$$

HOP

$$\pi_{\text{HOP}}(b) = \arg \max_{a_{t:\infty}^i} \frac{1}{m} \sum_{i=1}^m \gamma^+ R(s_t^i, a_t^i)$$

$$\text{subject to } s_{t+1}^i = G(s_t^i, a_t^i, w_i)$$

$$a_i^i = a_i^j \quad \forall i, j$$