

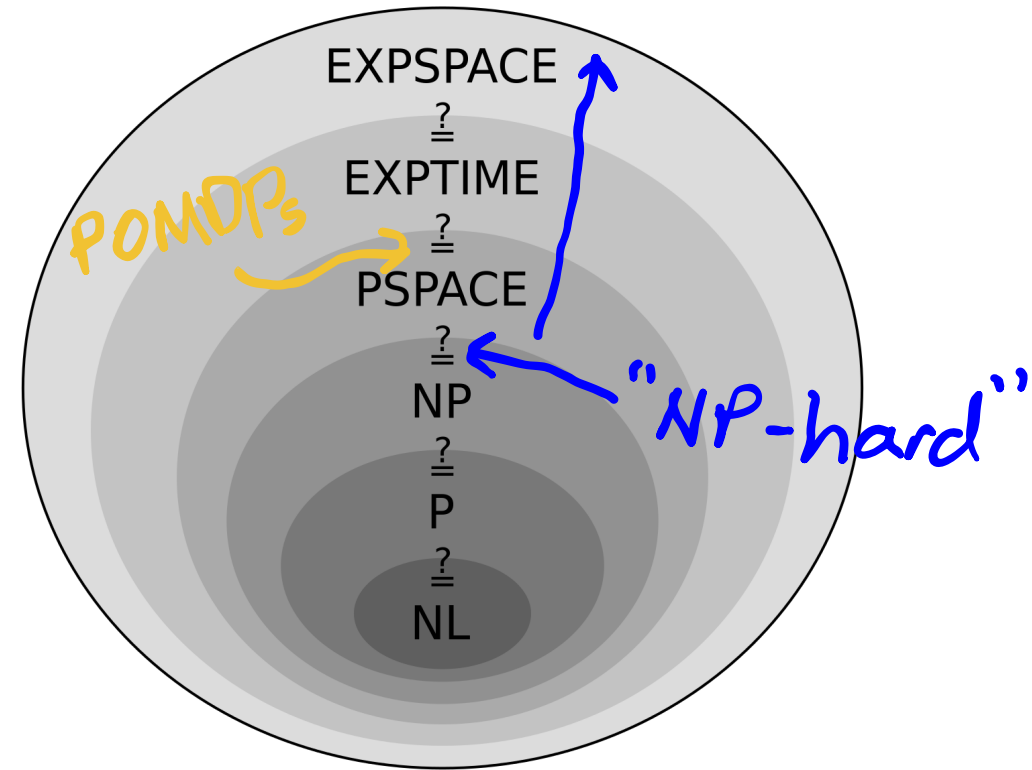
Recap

- Alpha Vectors
- Best solver for discrete POMDPs:

POMDP Computational Complexity

Sad 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 algorithm that can solve a general POMDP will have exponential complexity (we think)



Approximate POMDP Solutions

Numerical Approximations

(approximately solve original problem)



Offline

Tuesday



Online

After spring
break

Formulation Approximations

(solve a slightly different problem)

Today!

POMDP Objective

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

$$b' = \tau(b, a, o)$$

Certainty Equivalent

POMDP Objective

$$\pi(s) = -Ks$$

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

$$\pi_{\text{CE}}(b) = \pi_s(\mathbb{E}[s]_{s \sim b})$$

$$b' = \tau(b, a, o)$$

$$s_{t+1} = A s_t + B a_t + w_t$$

$$b' = \tau(b, a, o) \quad w_t \sim \mathcal{N}(0, \bar{W})$$

$$R = -s_t^T R_s s_t + -a_t^T R_a a_t$$

$$o_t = C s_t + v_t$$

$$v_t \sim \mathcal{N}(0, \bar{V})$$

Certainty Equivalent

Optimal for LQG

$$T(\mathbf{s}' \mid \mathbf{s}, \mathbf{a}) = \mathcal{N}(\mathbf{s}' \mid \mathbf{T}_s \mathbf{s} + \mathbf{T}_a \mathbf{a}, \Sigma_s)$$

$$O(\mathbf{o} \mid \mathbf{s}') = \mathcal{N}(\mathbf{o} \mid \mathbf{O}_s \mathbf{s}', \Sigma_o)$$

$$b(\mathbf{s}) = \mathcal{N}(\mathbf{s} \mid \boldsymbol{\mu}_b, \Sigma_b)$$

$$\begin{cases} \boldsymbol{\mu}_p \leftarrow \mathbf{T}_s \boldsymbol{\mu}_b + \mathbf{T}_a \mathbf{a} \\ \Sigma_p \leftarrow \mathbf{T}_s \Sigma_b \mathbf{T}_s^\top + \Sigma_s \end{cases}$$

$$\mathbf{K} \leftarrow \Sigma_p \mathbf{O}_s^\top \left(\mathbf{O}_s \Sigma_p \mathbf{O}_s^\top + \Sigma_o \right)^{-1}$$

$$\boldsymbol{\mu}_b \leftarrow \boldsymbol{\mu}_p + \mathbf{K}(\mathbf{o} - \mathbf{O}_s \boldsymbol{\mu}_p)$$

$$\Sigma_b \leftarrow (\mathbf{I} - \mathbf{K} \mathbf{O}_s) \Sigma_p$$

QMDP

POMDP Objective

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

$$b' = \tau(b, a, o)$$

Handwritten note: $\operatorname{argmax}_a Q_a^T b$

$$\pi_{\text{QMDP}}(b) = \operatorname{argmax}_{a \in A} \mathbb{E}_{s \sim b} [Q_{\text{MDP}}(s, a)]$$

$$b' = \tau(b, a, o)$$

Example: Tiger POMDP with Waiting

$$P(o = TL | a = \text{listen}, s = TL) = 0.85$$

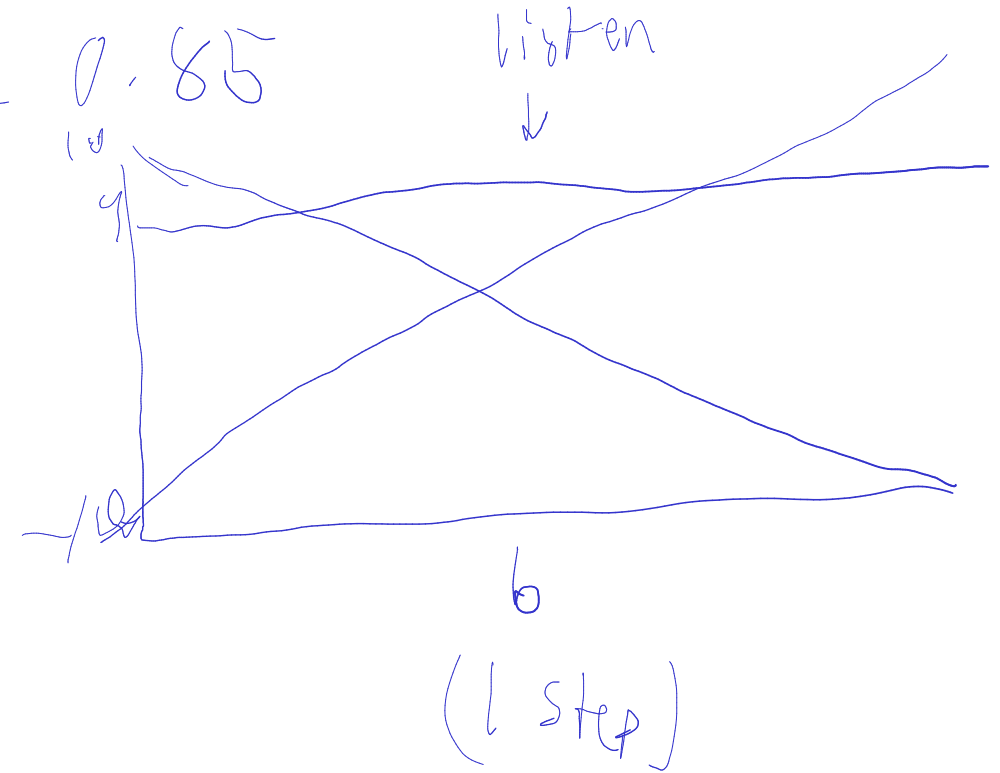
$$a = TL \quad s = TR$$

$$a = \text{listen}$$

$$a = TL \quad s = TL$$

$$a = \text{wait} \quad \leftarrow$$

$$R = \begin{cases} 10 \\ -1 \\ -100 \\ 0 \end{cases}$$

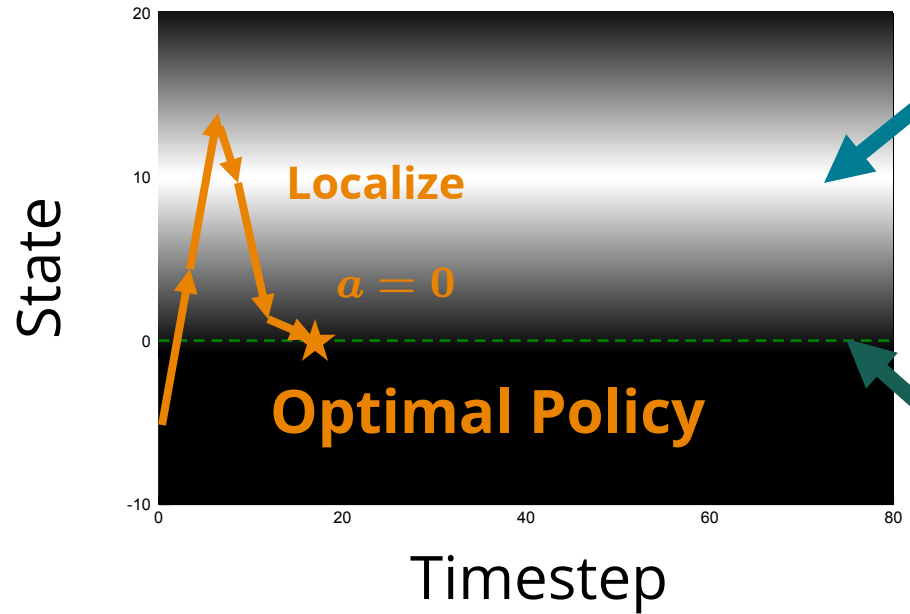


$$\gamma = 0.9$$

QMD Fails when information gathering is cast by

POMDP Example: Light-Dark

Accurate Observations

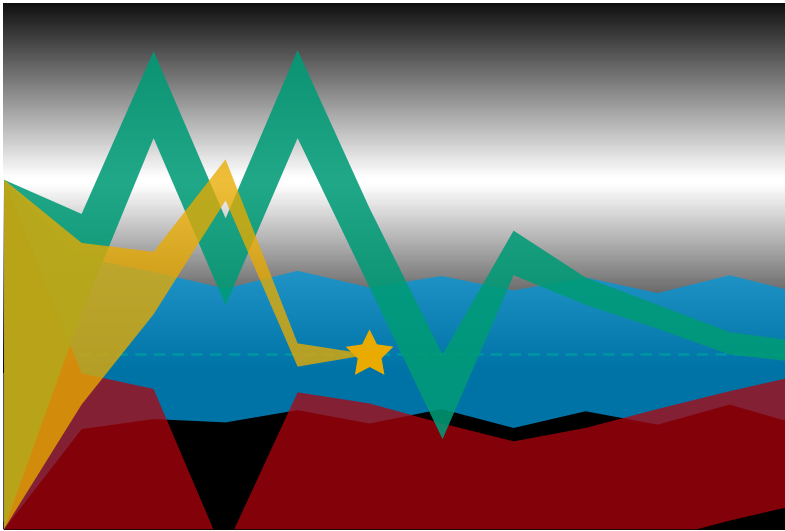


$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\}\end{aligned}$$

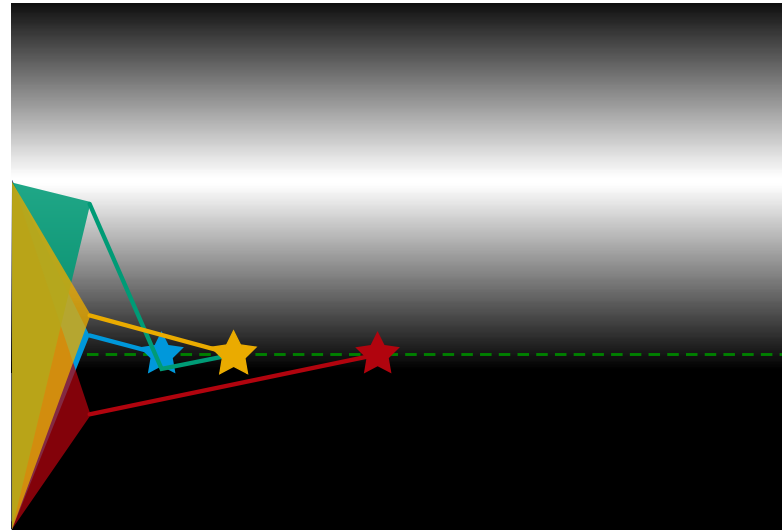
$$R(s, a) = \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}$$

Goal: $a = 0$ at $s = 0$

POMDP Solution



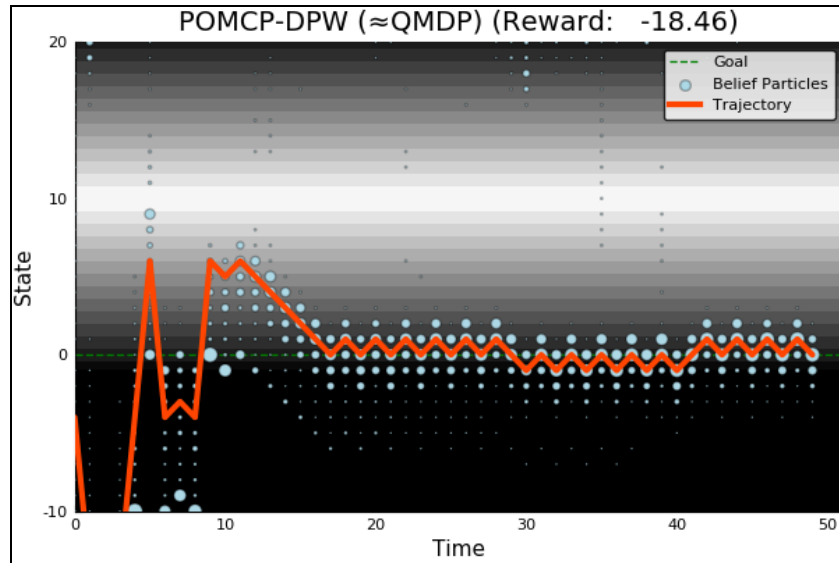
QMDP



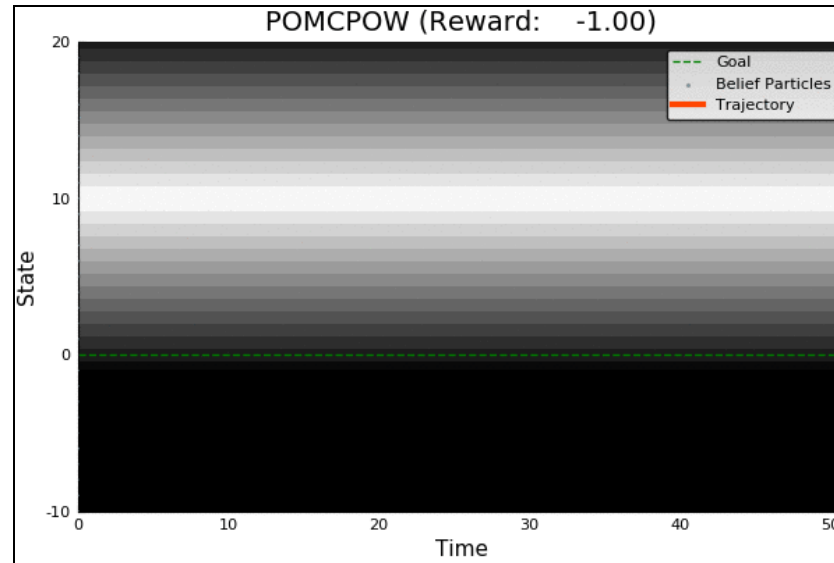
Same as **full observability**
on the next step

Information Gathering

QMDP

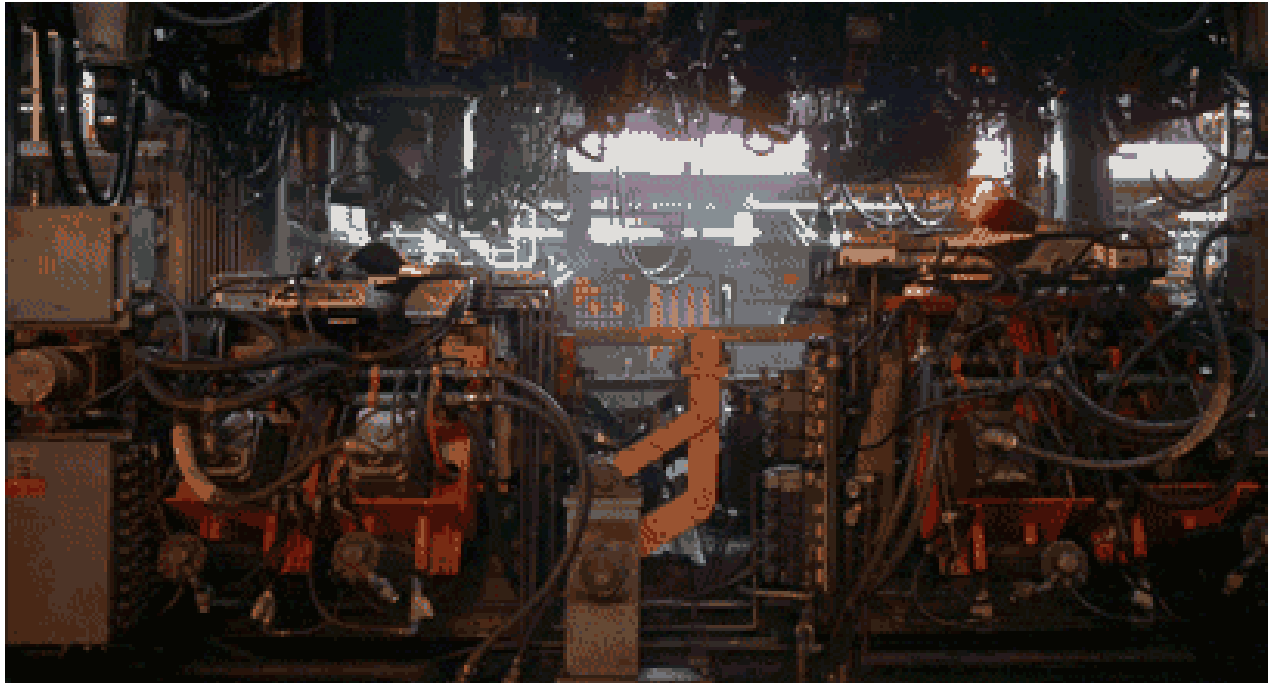


Full POMDP



QMDP

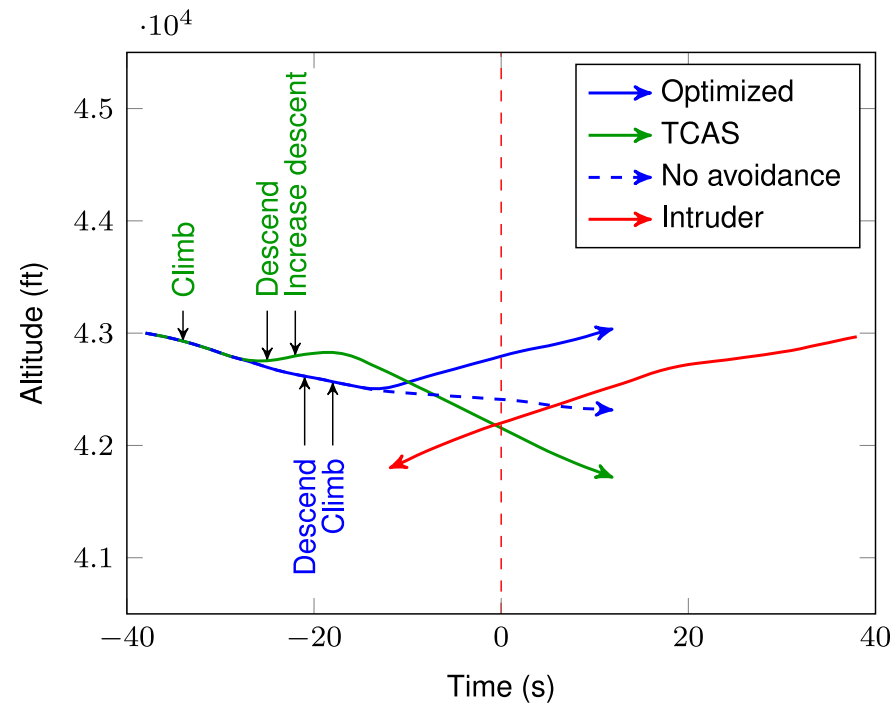
INDUSTRIAL GRADE



QMDP

ACAS X

[Kochenderfer, 2011]



Fast Informal Bound

POMDP Objective

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

$$b' = \tau(b, a, o)$$

QMDP (VI)

$$Q_a^{(k+1)} = R(s, a) + \gamma \sum_{s'} T(s' | s, a) \max_{a'} Q_{a'}^{(k)}[s']$$

FIB

$$Q_a^{(k+1)} = R(s, a) + \gamma \sum_{\theta} \max_{a'} \sum_{s'} T(s' | s, a, \theta) \cdot \underbrace{z(o | a, s')}_{\text{}} \cdot \underbrace{Q_{a'}^{(k)}[s']}_{\text{}}$$

Hindsight Optimization

POMDP Objective

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

$$b' = \tau(b, a, o)$$

$$s' = G(s, \underline{a}, w_t)$$

Choose w_t

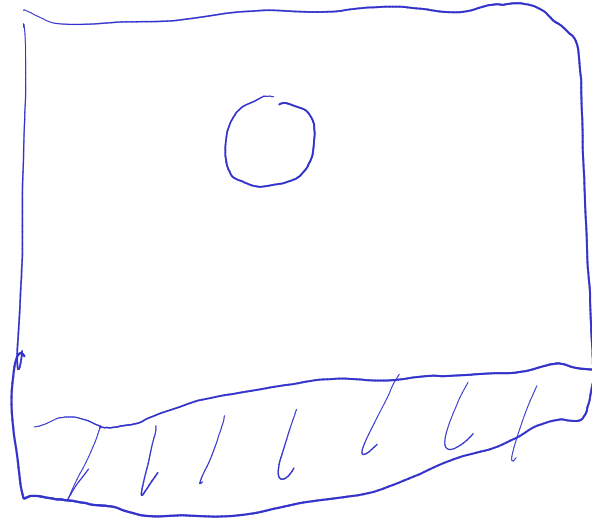
$$a = \operatorname{argmax}_{a_1, a_2, \dots} \mathbb{E} \left[R(s, a) + \dots \right]$$

k-Markov

POMDP Objective

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

$$b' = \tau(b, a, o)$$



$$s_t = [o_t, o_{t-1}, o_{t-k}]$$

Open Loop

POMDP Objective

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

$$b' = \tau(b, a, o)$$

$$(a_1, a_2, \dots) = \operatorname{argmax}_{(a_1, a_2, \dots)}$$

$$\mathbb{E} \left[\sum_{t=0}^T \gamma^t R(s_t, a_t) \right]$$

Comparison

Name	Description	Properties	Usefulness
Certainty Equivalence	Control as if the true state is mean of belief	Optimal for LQG	5 star
QMDP	Full observability after 1 time step	QMDP produces an upper bound for the true V	5 star
Hindsight Optimization	Hindsight knowledge of state and outcome uncertainty	Looser upper bound than QMDP	4 star
FIB	Takes 1 observation into account	Tighter upper bound than QMDP	2 star
k-Markov	Pretend that last k observations make up the state and solve that MDP	Great for Atari!	4 star
Open Loop	Choose sequence of actions	Good if aleatory is low, and epistemic is hard to reduce	3 star

Comparison

Name	Description	Properties	Usefulness
Certainty Equivalence	Control as if the true state is mean of belief	Optimal for LQG	5 star
QMDP	Full observability after 1 time step	QMDP produces an upper bound for the true V	5 star
Hindsight Optimization	Hindsight knowledge of state and outcome uncertainty	Looser upper bound than QMDP	4 star
FIB	Takes 1 observation into account	Tighter upper bound than QMDP	2 star
k-Markov	Pretend that last k observations make up the state and solve that MDP	Great for Atari!	4 star
Open Loop	Choose sequence of actions	Good if aleatory is low, and epistemic is hard to reduce	3 star