POMDP Model Approx.        "Formulation" Approx

Last two Lectures: Numerically Approximate Solutions
        — Offline    PBVI / $\alpha$ vectors
        — Online    tree search
                DESPOT $\alpha$

Model Approximation
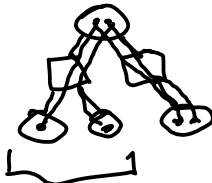        — Solve a slightly different opt. prob.
                                    "model"
        - Toolbox
        - when to use

## DESPOT-$\alpha$

POMCPOW: MCTS + Prog. Wid. + Weighted
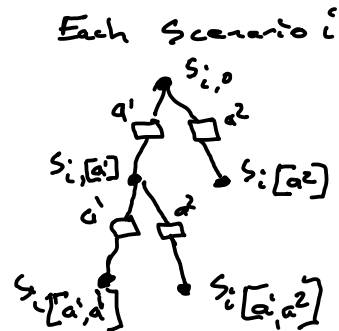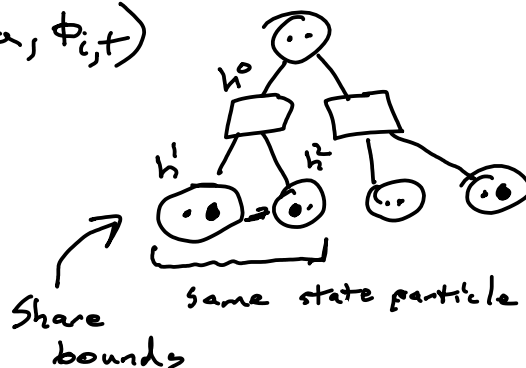                                    Particle Beliefs



DESPOT: Heuristic Search + Scenarios



## DESPOT-$\alpha$

All scenarios in *all* belief node
            weghted by $z$

$s', o, r = G(s, a, \phi_i, t)$



Share
bounds

Same state particle

Each Scenario $i$



$$U(h') = \sum_{i=1}^{K} w_i^{h'} U(s_{i, h^0})$$

$$L(h') = \sum_{i=1}^{K} w_i^{h'} \cdot \alpha^p [s_{i, h^0}]$$
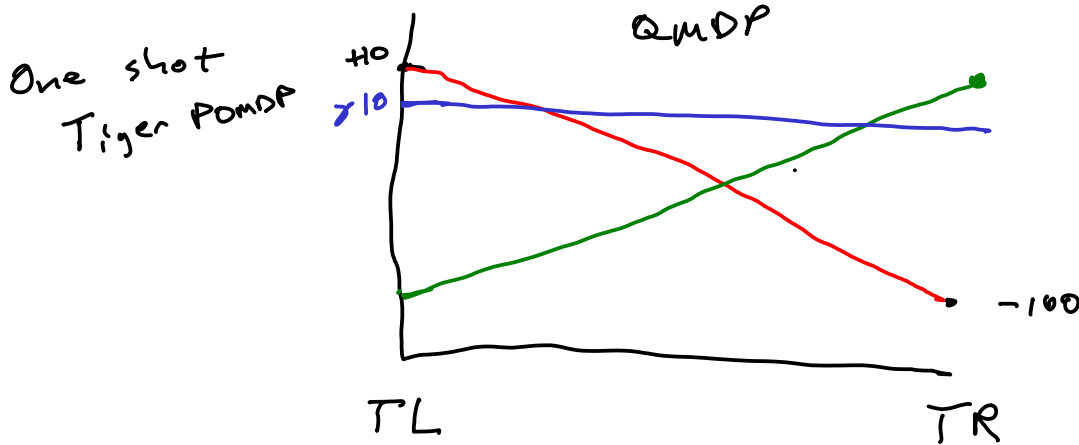
— any suboptimal
    plan

# Model Approximations

**QMDP**

$\alpha$-vectors $\sim$ Q values $\qquad \alpha^p[s] = Q^p(s,p[s])$

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

↖ solution to underlying MDP

QMDP

One shot
Tiger POMDP



+10

x10

−100

TL                                              TR

POMDP objective

$$\pi^* = \underset{\pi: B > A}{\arg\max} E\left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

QMDP

$$\pi_{QMDP}(b) = \underset{a \in A}{\arg\max} \underset{s \sim b}{E}\left[ Q_{MDP}(s,a) \right]$$

— belief at current time
↪ F.O. state in future

⌃ Optimistic

ACAS X          Fully Observable
                                    ↗ partially
$$s = (\Delta h, \dot{h}_0, \dot{h}_i, adv, ps)$$
                                    ↖ pilot state $\in \begin{Bmatrix} responsive \\ ignoring \end{Bmatrix}$
          └──── millions ────┘

QMDP
$$E_{s \sim b}\left[ Q_{MDP}(s,a) \right]$$

| Name | Desc. | Properties | Usefulness |
|---|---|---|---|
| QMDP | Full obs. after 1 step hindsight kn. of epistemic | Upper Bound on true Value Function | ★★★ ★★ ★ |
| FIB | takes 1 step observation into account | Tighter upper Bound than QMDP | ★ |
| Hindsight Optimization | Hindsight + knowledge of state + outcome (aleatory + epistemic) uncertainty | Looser upper bound than QMDP | ★★★★ |
| Certainty Equivalence | Control as if mean (or median or mode) is true state | Optimal for LQR | ★★★★★ |
| Open Loop action sequence - - - - - no observations | Choose action sequence that optimizes objective in expectation | Good when hard to reduce epistemic | ★★ ★ |
| Last k observations "k-markov" | Pretend last k observations are state and solve MDP | Great for Atari! | ★★★★ |
| Most Likely Observation | Plan assuming $b' = \tau(b, a, \hat{o}(b))$ | No observation branching Good when Z unimodal | ★★ ★ |
| Epistemic → Aleatory | Assume that partially observable part of s takes a random value at each time step | Conservative | ★★ |

# FIB

$$\pi_{FIB}(b) = \underset{a \in A}{argmax}\ \alpha_{FIB}^{a\ T} b$$

$$\alpha_a^{(k+1)}[s] = R(s,a) + \gamma \sum_o \underset{a'}{max} \sum_{s'} Z(o|a,s') T(s'|s,a) \alpha_a^{(k)}[s']$$

# Hindsight

### POMDP objective

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

### QMDP

$$\pi_{QMDP}(b) = \underset{a \in A_{sub}}{argmax} E\left[ Q_{MDP}(s,a) \right]$$

### HOP

$$V_{hs}(b) = E\left[ \underset{(a_1..a_T)}{max} \sum_{t=0}^{T} \gamma^t R(s_t, a_t) \right]$$

$$\pi_{hs}(b) = \underset{a \in A_{sub}}{argmax} E\left[ R(s_0,a) + \underset{\substack{a_t \\ t \in 1..T}}{max} \sum_{t=1}^{T} \gamma^t R(s_t, a_t) \right]$$

in VI
E

# CE

$$\pi_s: S \to A$$

$$\pi_{CE}(b) = \pi_s\left( \underset{sub}{E}[s] \right)$$

# O.L.

$$\underset{(a_1 ... a_T)}{max} E\left[ \sum_{t=0}^{T} \gamma^t R(s_t, a_t) \right]$$

no observation

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