

## Last Time

### POMDP Policies

└ History or Belief-Based

Approximate

Particle Filter

## This Time

### POMDP policy structure

└ Alpha Vectors

└ POMDP Value Iteration

$$\pi(b) = \begin{cases} OL & \text{if } b(TR) > 0.85 \\ OR & \text{if } b(TL) > 0.85 \\ L & \text{otherwise} \end{cases}$$

(DMU)

### Crying Baby Problem

$$S = \{h, \neg h\}$$

$$A = \{f, \neg f\}$$

$$O = \{c, \neg c\}$$

$$R(s, a) = -5 \mathbb{1}(a=f) - 10 \mathbb{1}(s=h)$$

$$T(h | \neg h, \neg f) = 0.1$$

$$T(\neg h | \cdot, f) = 1.0$$

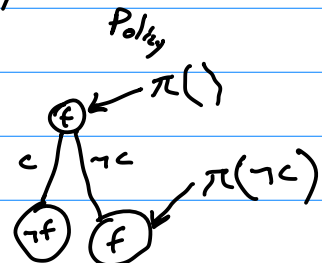
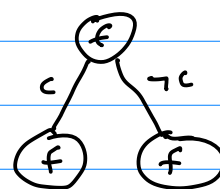
$$Z(c | \cdot, \neg h) = 0.1$$

$$Z(c | \cdot, h) = 0.8$$

Conditional Plan: Fixed-Depth History Based  
Fixed Initial Action

1 step:  $(f)$   $(\neg f)$

2 step:



# of possible conditional plans

$$|A|^{(101^h - 1) / (101 - 1)}$$

$$2^3 = 8$$

Value of a conditional plan  $\pi$  at state  $s$   
 $V^\pi(s)$

1 step

$$V^\pi(s) = R(s, \pi(s))$$

$$\pi = \textcircled{f}$$

$$V^\pi(h) = -15$$

$$V^\pi(\neg h) = -5$$

$$\pi = \textcircled{f}$$

$$V^\pi(h) = -10$$

$$V^\pi(\neg h) = 0$$

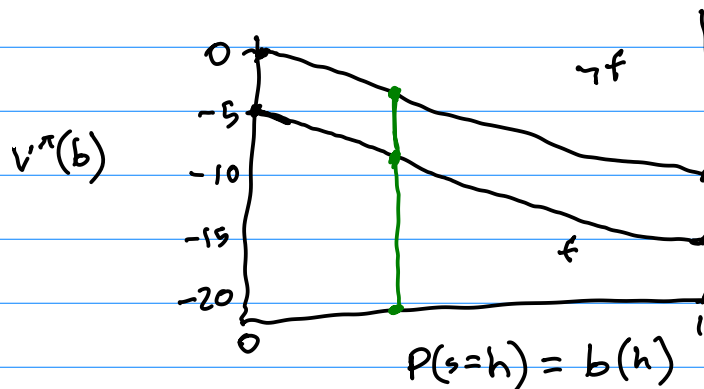
belief  $b$

$$V^\pi(b) = \sum_s V^\pi(s) b(s) = E_{s \sim b} [V^\pi(s)]$$

Alpha Vector

$$\alpha_\pi[s] = V^\pi(s)$$

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



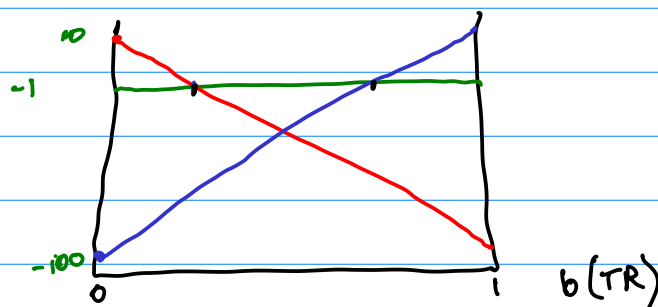
$$b(h) = 0.35$$

Tiger POMDP

$$S = \{TL, TR\}$$

$$A = \{OL, OR, L\}$$

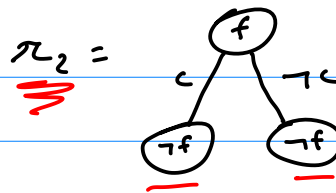
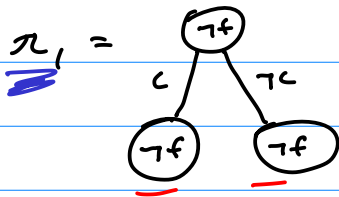
$$R = \begin{cases} +10 & \text{if open non-tiger door} \\ -100 & \text{if open tiger door} \\ -1 & \text{if listen} \end{cases}$$



$\textcircled{OL}$

$\textcircled{OR}$

$\textcircled{L}$



→  $V^\pi(s) = R(s, \pi()) + \gamma \left[ \sum_{s'} T(s'|s, \pi()) \sum_{\phi} \cancel{z(\phi|\pi(), s')} \cancel{V^{\pi(\phi)}(s')} \right]$   $R(s', \gamma f)$

$$V^\pi(s) = R(s, \pi()) + \gamma \sum_{s'} T(s'|s, \pi()) R(s', \gamma f)$$

$$V^{\pi_1}(h) = -10 + 0.9(1.0 \cdot -10) = -19$$

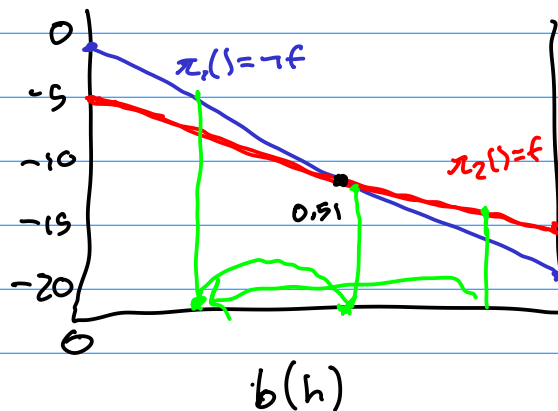
$$V^{\pi_1}(\neg h) = 0 + 0.9(0.9 \cdot 0 + 0.1 \cdot -10) = -0.9$$

$$V^{\pi_2}(h) = -15 + 0.9(1.0 \cdot 0) = -15$$

$$V^{\pi_2}(\neg h) = -5 + 0.9(1.0 \cdot 0) = -5$$

$$\alpha_{\pi_1} = [-0.9, -19]$$

$$\alpha_{\pi_2} = [-5, -15]$$



$$\pi(b) = \left[ \underset{\pi_i}{\operatorname{argmax}} \alpha_i^T b \right](\cdot)$$

In general

$$V^*(b) = \max_{\alpha \in \Gamma^*} \alpha^T b$$

↑ set of  $\alpha$  vectors

Convex, Piecewise  
Linear  
Value Functions

# POMDP Value Iteration

Starting with 1 step plans

Creating all Possible 2-step Plans

Remove Suboptimal Plans

Create 3 step plans

⋮

loop

$$\Gamma^l = \bigcup_{a \in A} \Gamma^a$$

$$\text{where } \Gamma^a = \bigoplus_{o \in O} \Gamma^{a,o}$$

$$\text{where } \Gamma^{a,o} = \left\{ \frac{1}{|O|} r_a + \gamma \alpha^{a,o} : \alpha \in \Gamma \right\}$$

$$\text{where } \alpha^{a,o}[s] = \sum_{s'} \mathbb{E}(o|a, s') T(s'|s, a) \alpha[s']$$

MDP Value Iteration

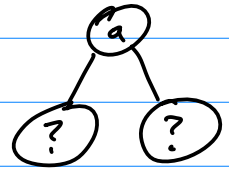
loop

$$V' \leftarrow B[V]$$

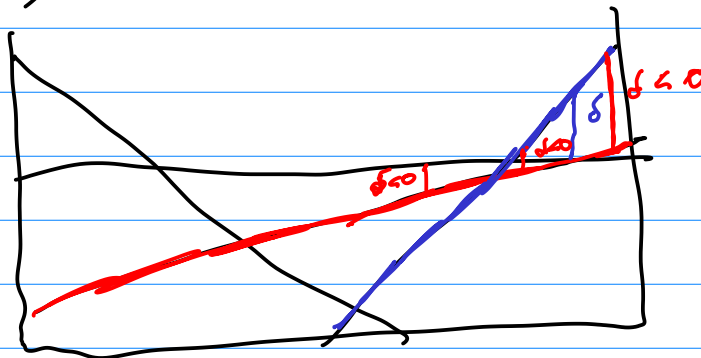
$$B[V](s) = \max_a R(s, a) + \gamma E[V(s')] \\ s' \sim T(s, a)$$

$$\Gamma_1 \oplus \Gamma_2 = \{ \alpha_1 + \alpha_2 : \alpha_1 \in \Gamma_1, \alpha_2 \in \Gamma_2 \}$$

$$r_a[s] = R(s, a)$$



Pruning



$$\begin{aligned} &\text{maximize } \delta \\ &\text{subject to } b \geq 0 \\ &\quad 1^T b = 1 \end{aligned}$$

$$\alpha^T b \geq \alpha'^T b + \delta \quad \forall \alpha' \in \Gamma$$

if  $\delta > 0$ , then  $\alpha$  is not dominated  
if  $\delta \leq 0$ ,  $\alpha$  is dominated, we can prune  
 $b$  that maximizes  $\delta$  is called "Witness  
belief"

- WITNESS
- Incremental Pruning