<u>Last Time</u>

$\alpha$ vectors

$\alpha_x[s] \cdot [0.4, 0.3, 0.2]$

set $\Gamma$

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

$S = \{1, 2\}$



$\max_{\alpha \in \Gamma} \alpha^T [0.1, 0.9]$

$\alpha$   $b$

0   $b(2)$   1

<u>Today</u>

Offline POMDP

"Survey of Point · Based POMDP Solvers"

| Name | S | A | O | |
|------|---|---|---|---|
| Exact VI | D (10) | D | D | |
| PBVI | D | D | D | bigger |
| Perseus | D | D | D | |
| HSVI | D | D | D | |
| SARSOP | D (1000) | D | D | |
| MCVI | C | D | D/C | Linear (Gaussian Quadratic) |
| LQG | C | C | C | |

Exact VI

loop

$$\Gamma' \doteq \bigcup_{a \in A} \Gamma^a$$

$$\Gamma^a = \bigoplus_{o \in O} \Gamma^{ao}$$

for 1 iteration

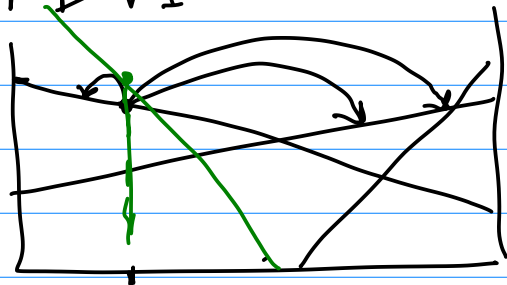$$O\left( |\Gamma||A||O||S|^2 + |A||S||\Gamma|^{|O|}\right)$$

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

$$\alpha^{a,o}[s] = \sum_{s'} Z(o|a, s') \, T(s'|s, a) \, \alpha[s']$$

# PBVI



$\rightarrow$ Backup Belief $(\Gamma, b)$

    for $a \in A$

        for $o \in O$

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

$$\alpha_{a,o} = \operatorname*{argmax}_{\alpha \in \Gamma} \alpha^T b'$$

        for $s \in S$

$$\alpha_a[s] = R(s,a) + \gamma \sum_{s',o} T(s'|s,a) \, Z(o'|a,s') \, \alpha_{a,o}[s']$$

return $\operatorname*{argmax}_{\alpha_a} \alpha_a^T b$

$$O\left(|\Gamma||A||O||S|^2 + |A||S||\Gamma||B|\right)$$

How we choose $B$

Original PBVI

$$B = \{b_o\}$$

loop $B' = \phi$

    for $b \in B$
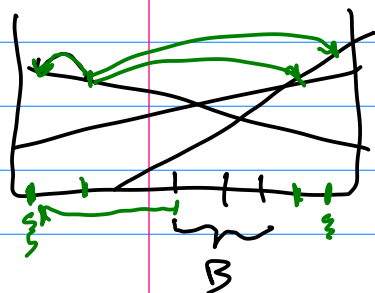
        $\Gamma = \Gamma \cup$ Backup Belief $(\Gamma, b)$

    for $b \in B$

        $\tilde{B} = \{\tau(b,a,o) : a \in A, o \in O\}$

        $B' = B' \cup \operatorname*{argmax}_{b \in \tilde{B}} \|B, b'\|$     $\leftarrow$ add furthest b away from B

    $B = B \cup B'$



$B = \{b^1, b^2, b^3\}$

$B' = \{b^4, b^5\}$

$B \cup B' = \{b^1 \dots b^5\}$

Perseus — randomly choose B

→ . $\quad B = \emptyset$

→ $\quad b = b_0$

loop until $|B| = n$

$\qquad a = \text{rand}(A)$

$\qquad o = \text{rand}(P(o|b,a))$

$\qquad B = B \cup \{\tau(b,a,o)\}$

---

## HSVI — Heuristic Search Value Iteration

$\bar{V}(b)$ . upper bound

$\underline{V}(b)$ lower bound

while $\widehat{V}(b_0) - \underline{V}(b_0) > \varepsilon$

$\qquad$ explore $(b_0, 0)$

explore $(b, t)$

$\quad$ if $\bar{V}(b) - \underline{V}(b) > \varepsilon \gamma^T$

$\qquad a^* = \underset{a}{\text{argmax}} \; \bar{Q}(b,a)$ ← upper bound

$\qquad$ excess uncertainty

$\qquad o^* = \underset{o}{\text{argmax}} \left( P(o|b,a) \left( \widehat{V}(\tau(b,a,o)) - \underline{V}(\tau(b,a,o) - \varepsilon \gamma^{(t+1)}) \right) \right.$
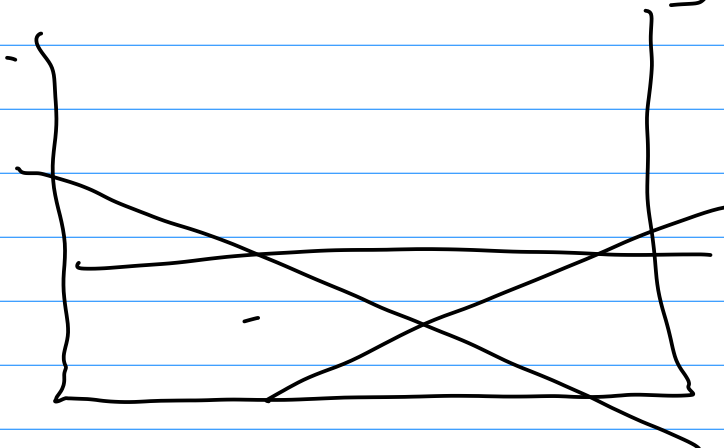
$\qquad$ explore $(\tau(b, a^*, o^*), t+1)$

→

⟶ $\qquad \underline{\Gamma} = \underline{\Gamma} \cup \text{Backup Belief}(\underline{\Gamma}, b)$

$\qquad \overline{V}(b) = B_b[\bar{V}(b)]$
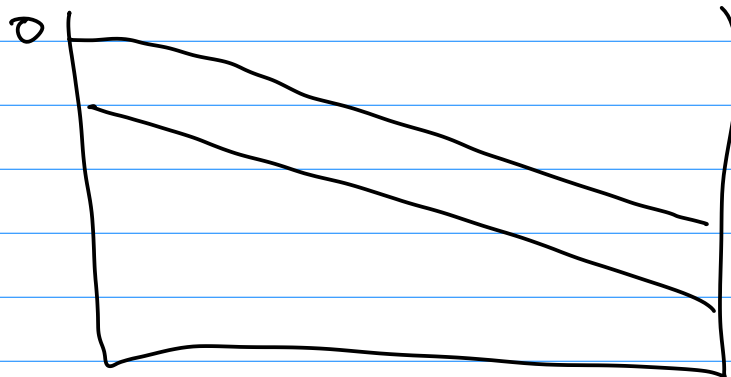
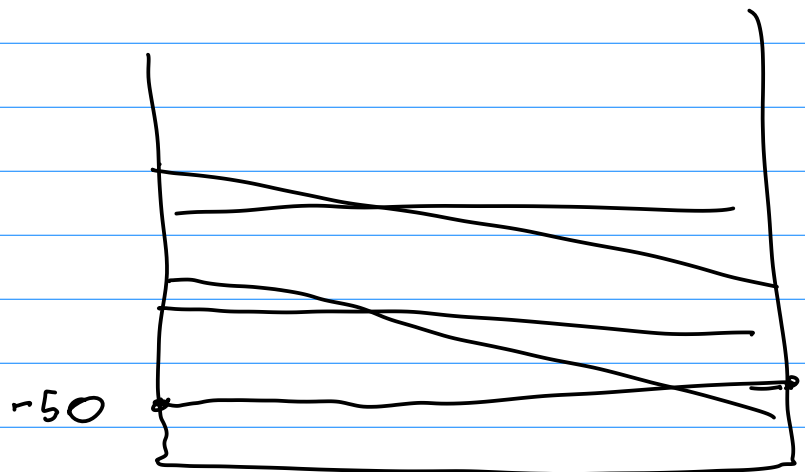Baby POMDP                    ⓕ          ⓣ̄ꜰ



$\dfrac{R}{1-\gamma}$          Open-Loop Policy

                              Never Feed

$\dfrac{-5}{1-0.9} = -50$



$-50$

Upper Bounds                                    Finding points
                                                to interpolate
        Sawtooth Upperbound                     between is
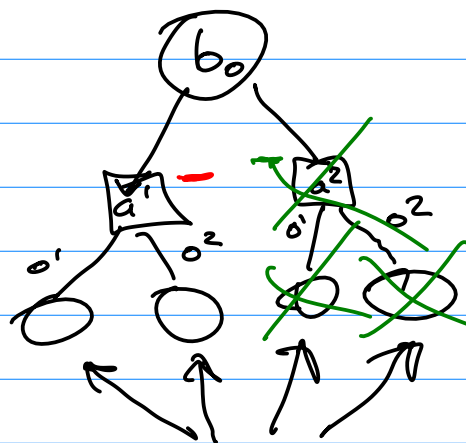                                                an LP

$$B_b[\bar{V}(b)] = \max_a R(b,a) + \gamma \sum_o P(o|b,a)\bar{V}(\tau(b,a,o))$$

SARSOP

Successive Approximation of Reachable Space under Optimal Policies
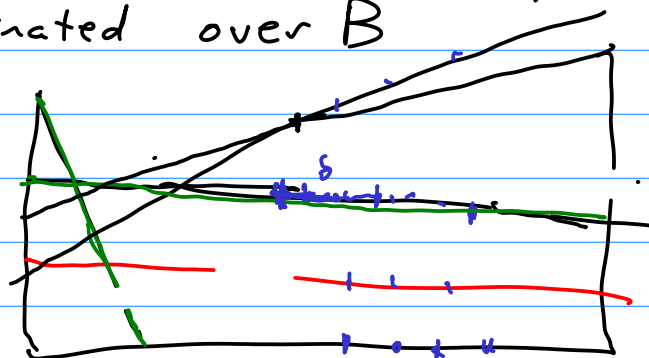
Similar to HSVI

$B \subset R \qquad B \subset R^*$



Belief points that could be in B

if $\bar{Q}(b,a^1) < \underline{Q}(b,a^2)$

then prune all $b$ below $(b,a^2)$ from B

Instead of pruning $\alpha$ that are dominated over the whole belief space, prune $\alpha$ dominated over B
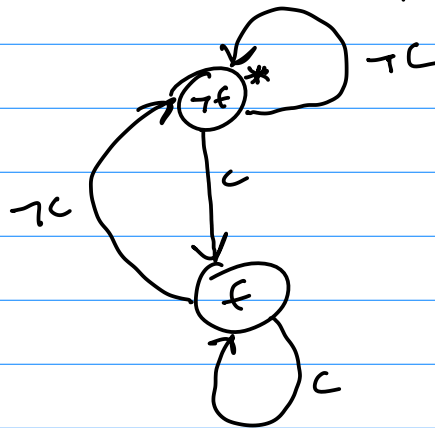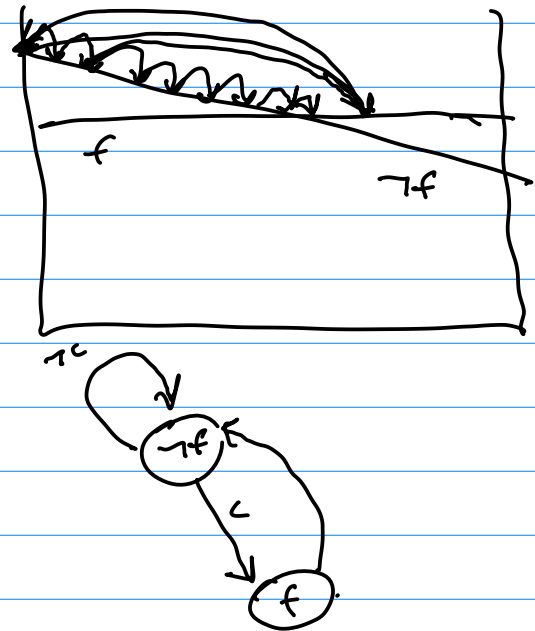


Pruned by Any

Pruned by SARSOP

B

# Policy Graph

vertex labeled with a
edge labeled with o

## Feed when Crying



Policy Graph → α vector
evaluate plan
α vectors → policy graphs



# MCVI   "Monte Carlo Value Iteration"

$$V_G(b) = \max \int_{s \in S} \alpha_v(s) \, b(s) \, ds \qquad \leftarrow \text{Approx w/MC}$$

$\vee$

MC-Backup $(G, b, N)$

$\quad R_a = 0 \quad V_{a,o,v} = 0$

$\quad$ for $a \in A$

$\quad\quad$ for $i$ in $1:N$

$\quad\quad\quad s_i \Leftarrow sample(b)$

$\quad\quad\quad s'_i, o_i, r_i \Leftarrow G(s_i, a)$

$\quad\quad\quad R_a + r_i$

$\quad\quad\quad$ for $v \in G$

$\quad\quad\quad\quad V_{a, o_i, v} = V_{a, o_i, v} + Simulate(G, v, s', L)$

$\quad\quad$ for $o$ in $O$

$\quad\quad\quad V_{a,o} = \max_{v \in G} V_{a,o,v}$

$\quad\quad\quad V_{a,o} = \operatorname{argmax}_{v \in G} V_{a,o,v}$

$\quad\quad V_a = R_a + \gamma \sum_o V_{ao}/N$

$V^* = \max_a V_a$

$a^* = \operatorname{argmax} V_a$

add new node to $G$ labeled with $a^*$



$V = value$

$v = $ node in graph