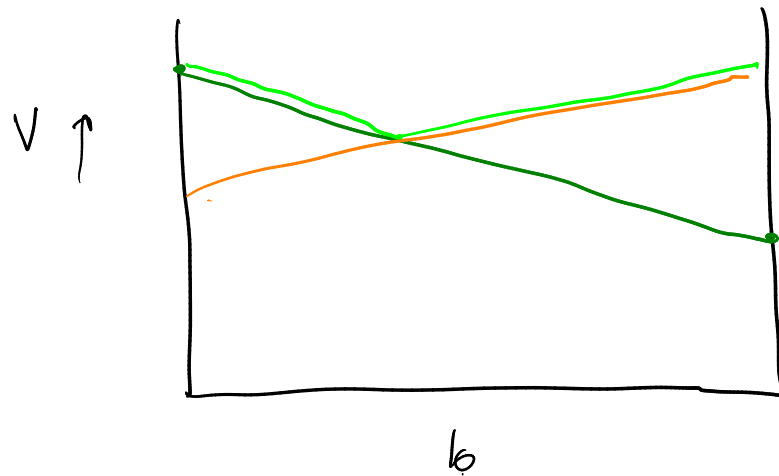# Offline POMDP Algorithms

# Last time: POMDP Value Iteration (horizon $d$)

$\Gamma^0 \leftarrow \emptyset$
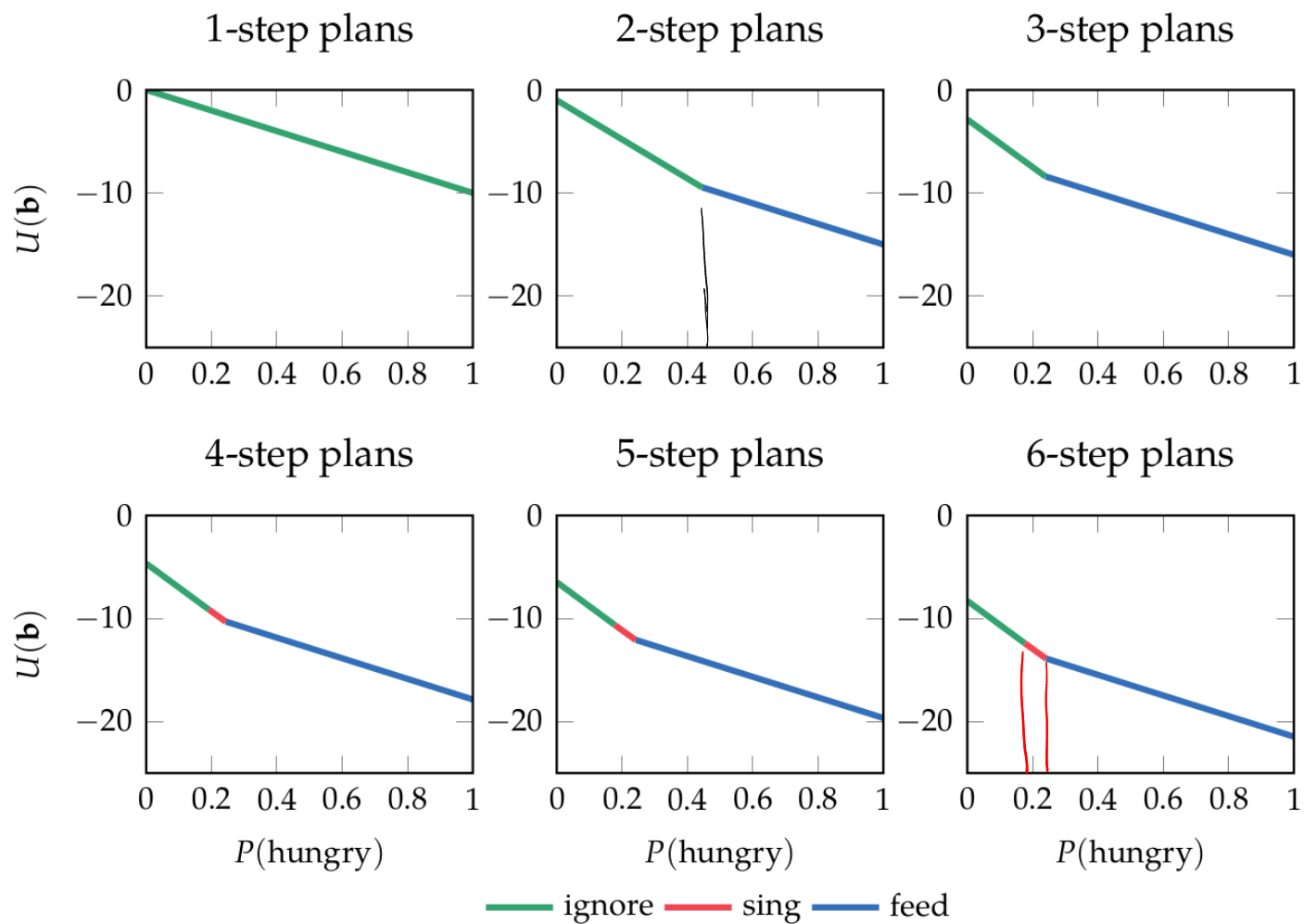
$|S| \approx 10$

for $n \in 1 \dots d$
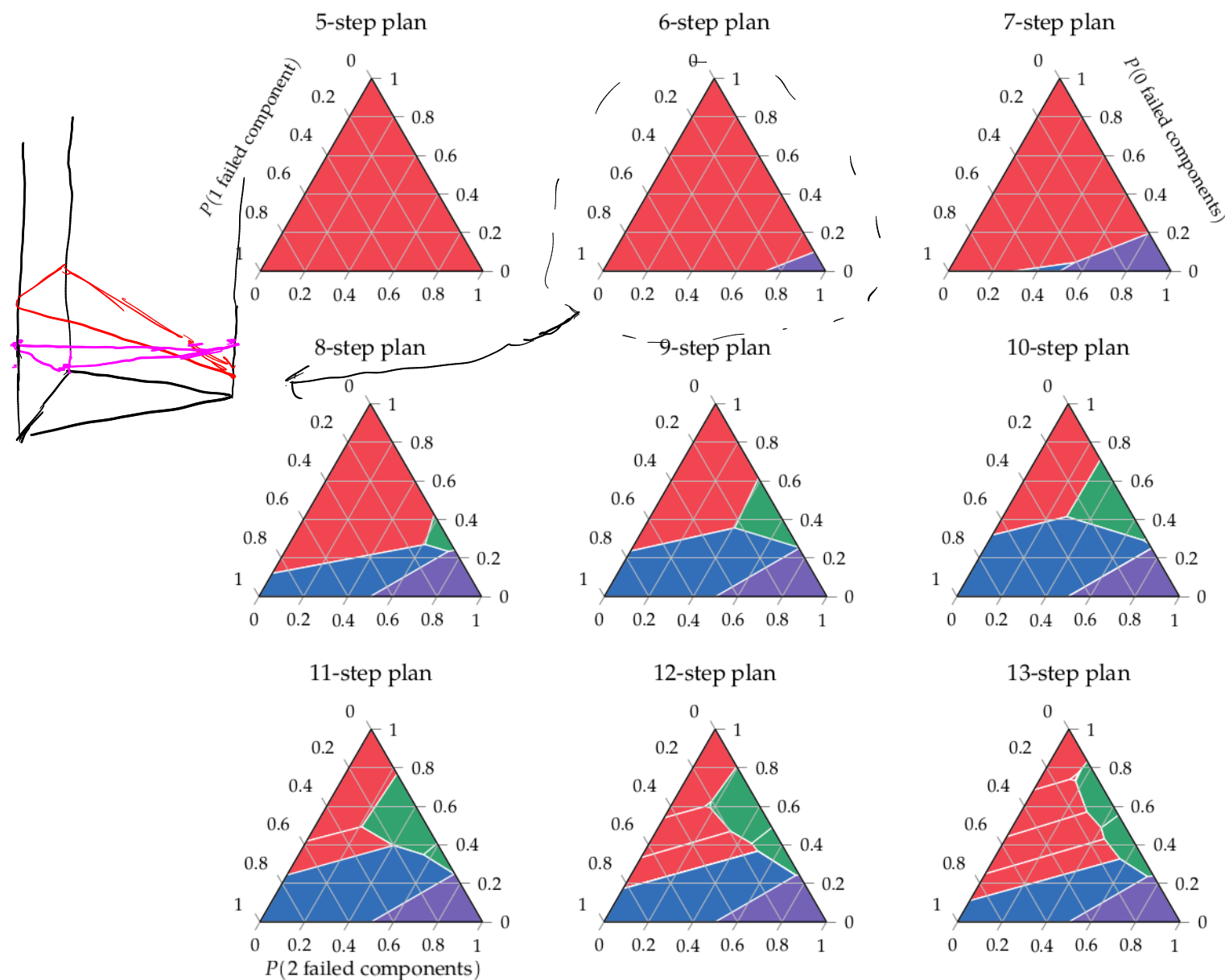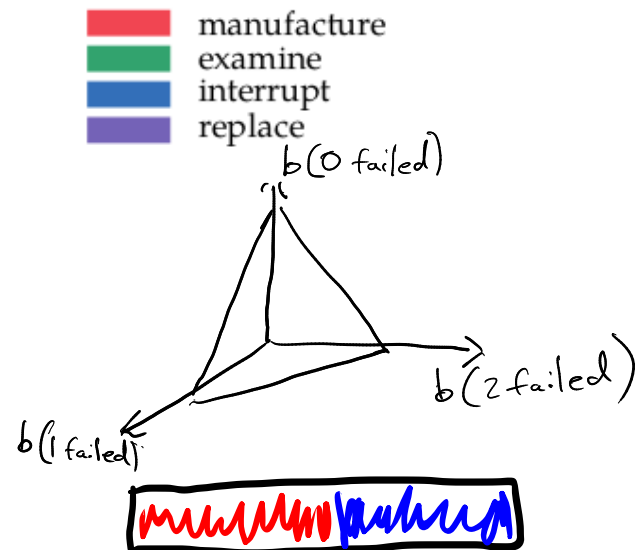
    Construct $\Gamma^n$ by expanding with $\Gamma^{n-1}$

    Prune $\Gamma^n$

# Finite Horizon POMDP Value Iteration

# Finite Horizon POMDP Value Iteration



5-step plan · 6-step plan · 7-step plan · 8-step plan · 9-step plan · 10-step plan · 11-step plan · 12-step plan · 13-step plan

$P(1 \text{ failed component})$ · $P(0 \text{ failed components})$ · $P(2 \text{ failed components})$

Legend:
- manufacture
- examine
- interrupt
- replace

$|S| = 3$

$b(0 \text{ failed})$

$b(1 \text{ failed})$

$b(2 \text{ failed})$

# Infinite-Horizon POMDP Lower Bound Improvement

.

# Infinite-Horizon POMDP Lower Bound Improvement

$$\alpha_a = (I - \gamma T^a)^{-1} R^a$$

always execute same action

$\Gamma \leftarrow$ blind lower bound

loop

$\quad \Gamma \leftarrow \Gamma \cup \text{backup}(\Gamma)$

$\quad \Gamma \leftarrow \text{prune}(\Gamma)$

A survey of point-based POMDP solvers

backup

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

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

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

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

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

$$\Gamma' \oplus \Gamma^2 = \{ \alpha_1 + \alpha_2 : \alpha_1 \in \Gamma', \alpha_2 \in \Gamma^2 \}$$

# Point-Based Value Iteration (PBVI)

point_backup$(\Gamma, b)$
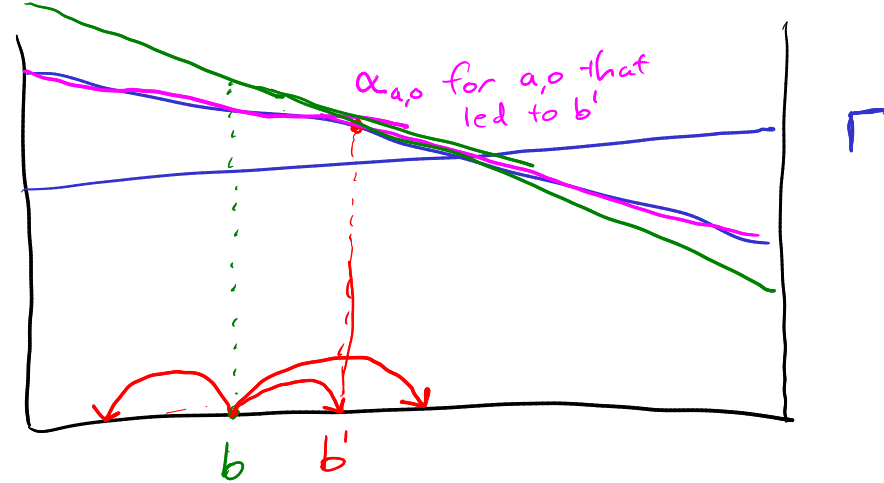
for $a \in A$

for $o \in O$

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

$\alpha_{a,o} \leftarrow \operatorname*{argmax}_{\alpha \in \Gamma} \alpha^\top b'$

for $s \in S$

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

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



$\alpha_{a,o}$ for $a,o$ that led to $b'$

$\Gamma$

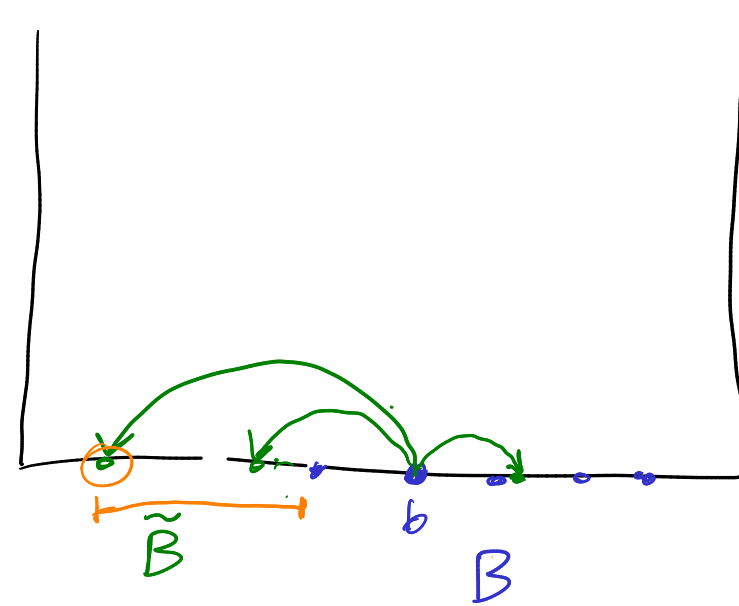$b$  $b'$

# Original PBVI

How do we choose B

$B \leftarrow b_0$

loop

   for $b \in B$

      $\Gamma \leftarrow \Gamma \cup \{\text{point\_backup}(\Gamma, b)\}$

   $B' \leftarrow \emptyset$

   for $b \in B$

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

      $B' \leftarrow B' \cup \left\{ \underset{b' \in \tilde{B}}{\text{argmax}} \, \|B, b'\| \right\}$

   $B \leftarrow B \cup B'$

# PERSEUS: Randomly Selected Beliefs

Two Phases:

    1. Random Exploration
    2. Value Backup

Random Exploration:

$B \leftarrow \emptyset$

$b \leftarrow b_0$

loop until $|B| = n$

    $a \leftarrow \mathrm{rand}(A)$

    $o \leftarrow \mathrm{rand}(P(o \mid b, a))$

    $b \leftarrow \tau(b, a, o)$

    $B = B \cup \{b\}$

# Heuristic Search Value Iteration (HSVI)

$\underline{V}(b_0)$ $\qquad$ $\overline{V}(b_0)$

while $\overline{V}(b_0) - \underline{V}(b_0) > \epsilon$

explore$(b_0, 0)$

explore(b, t)

  if $\overline{V}(b) - \underline{V}(b) > \epsilon\gamma^t$

promising

    $a^* = \underset{a}{\operatorname{argmax}}\ \overline{Q}(b, a)$

observations
1. likely
2. lots of uncertainty

    $o^* = \underset{o}{\operatorname{argmax}}\ P(o \mid b, a^*)\left(\overline{V}(\tau(b, a^*, o)) - \underline{V}(\tau(b, a^*, o)) - \epsilon\gamma^t\right)$

Weighted
Excess Uncertainty

    explore$(\tau(b, a^*, o^*), t+1)$

    $\underline{\Gamma} \leftarrow \underline{\Gamma} \cup \operatorname{point\_backup}(\underline{\Gamma}, b)$

    $\overline{V}(b) = B_b\left[\overline{V}(b)\right]$

$\Gamma = \{\alpha_1, \alpha_2 \ldots \alpha_n\}$

$B = \{b_1, b_2 \ldots b_m\}$ $\qquad \overline{V} =$ $\qquad \underline{V} =$

$\Gamma$

B

$\Gamma$ $\qquad$ V

# How do we represent an upper bound with alpha vectors
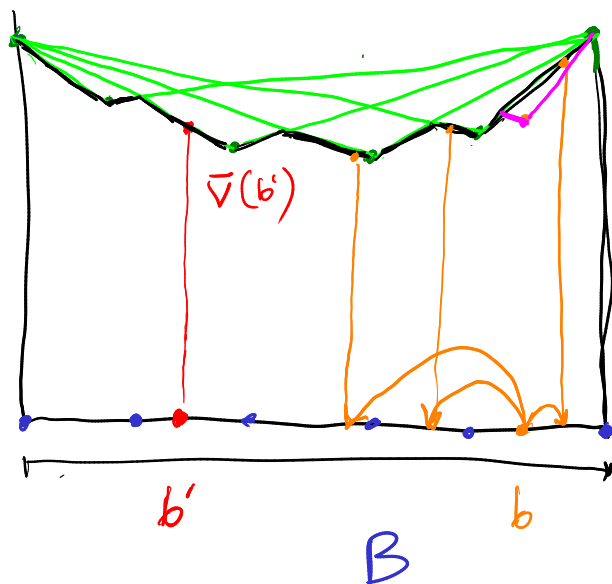


$$V^* \quad \overline{V}(b) = \min_{\alpha \in \widetilde{\Gamma}} \alpha^\top b$$

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

# Sawtooth Upper Bounds

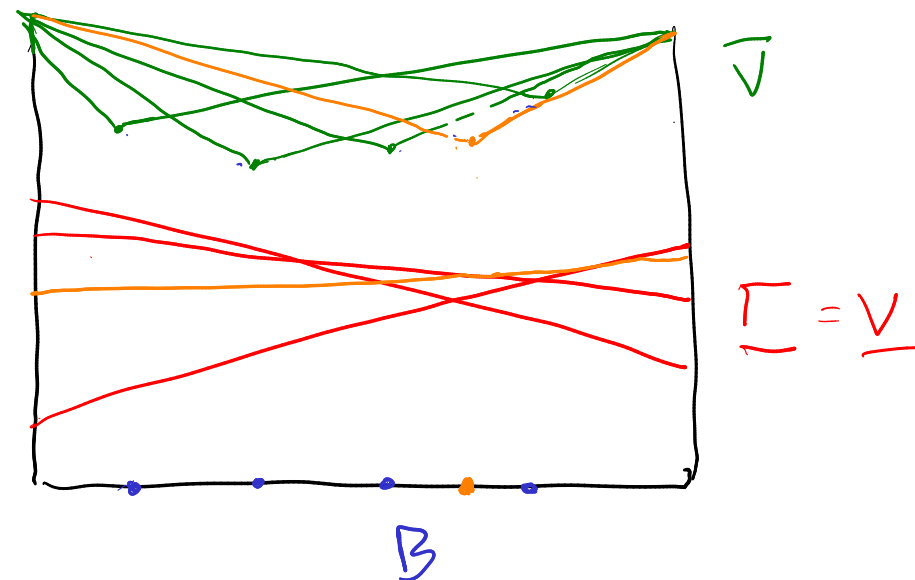for each $b \in \underline{B}$, store $\overline{V}(b)$

and the vertices of belief simplex

$\overline{V}(b)$ for $b \in B \cup$ vertices

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

$\overline{V}(b')$

$b'$

$b$

$B$

$\dfrac{\overline{R}}{1-\gamma} = \max_{a,s} R(s,a)$

What HSVI is working with

$\overline{V}$

$\underline{\Gamma} = \underline{V}$

$B$

# SARSOP

## Successive Approximation of Reachable Space under Optimal Policies

Similar to HSVI

HSVI

$B \subset R$

↳ reachable

SARSOP

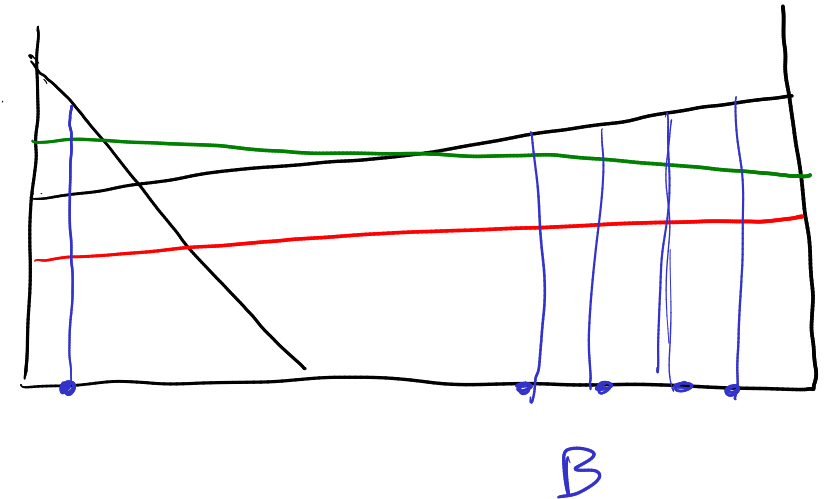$B \subset R^*$

↑ beliefs
reachable
under optimal
policy

+ a few other tricks

Witness
would prune
both //



$\overline{Q}(b,a)$     $\underline{Q}(b,a)$

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

then remove all $b$ below $(b,a^1)$ from $B$

B

Witness: ~20 states

SARSP: 10,000 − 100,000
states

# Offline POMDP Algorithms

# Policy Graphs

# Monte Carlo Value Iteration (MCVI)

$$\text{MC-Backup}\ (G, b, N)$$

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

for $a \in A$

    for $i$ in $1:N$

        $s_i \Leftarrow \text{sample}\ (b)$

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

        $R_a + r_i$

        for $v \in G$

            $V_{a,o_i,v} = V_{a,o_i,v} + \text{Simulate}\left(G, v, s'_i, L\right)$

    for $o$ in $O$

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

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

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

$$V^* = \max_a V_a$$

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

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