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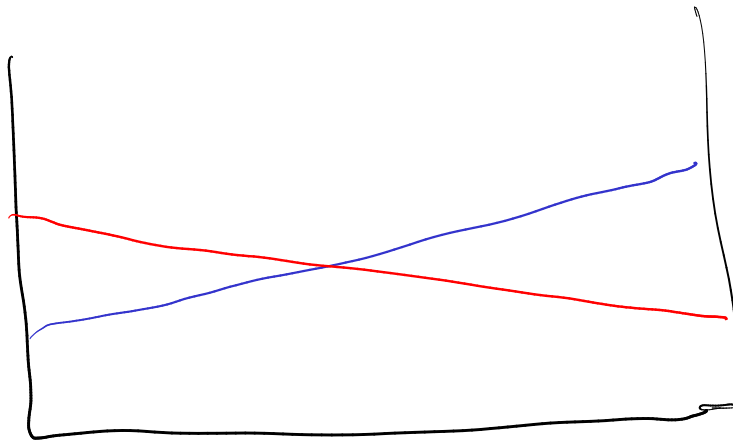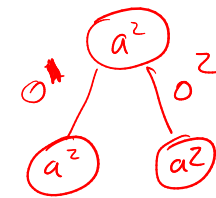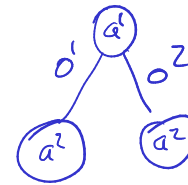
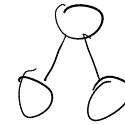# Last time: POMDP Value Iteration (horizon $d$)

$\Gamma^0 \leftarrow \emptyset$

for $n \in 1 \ldots d$

    Construct $\Gamma^n$ by expanding with $\Gamma^{n-1}$

    Prune $\Gamma^n$
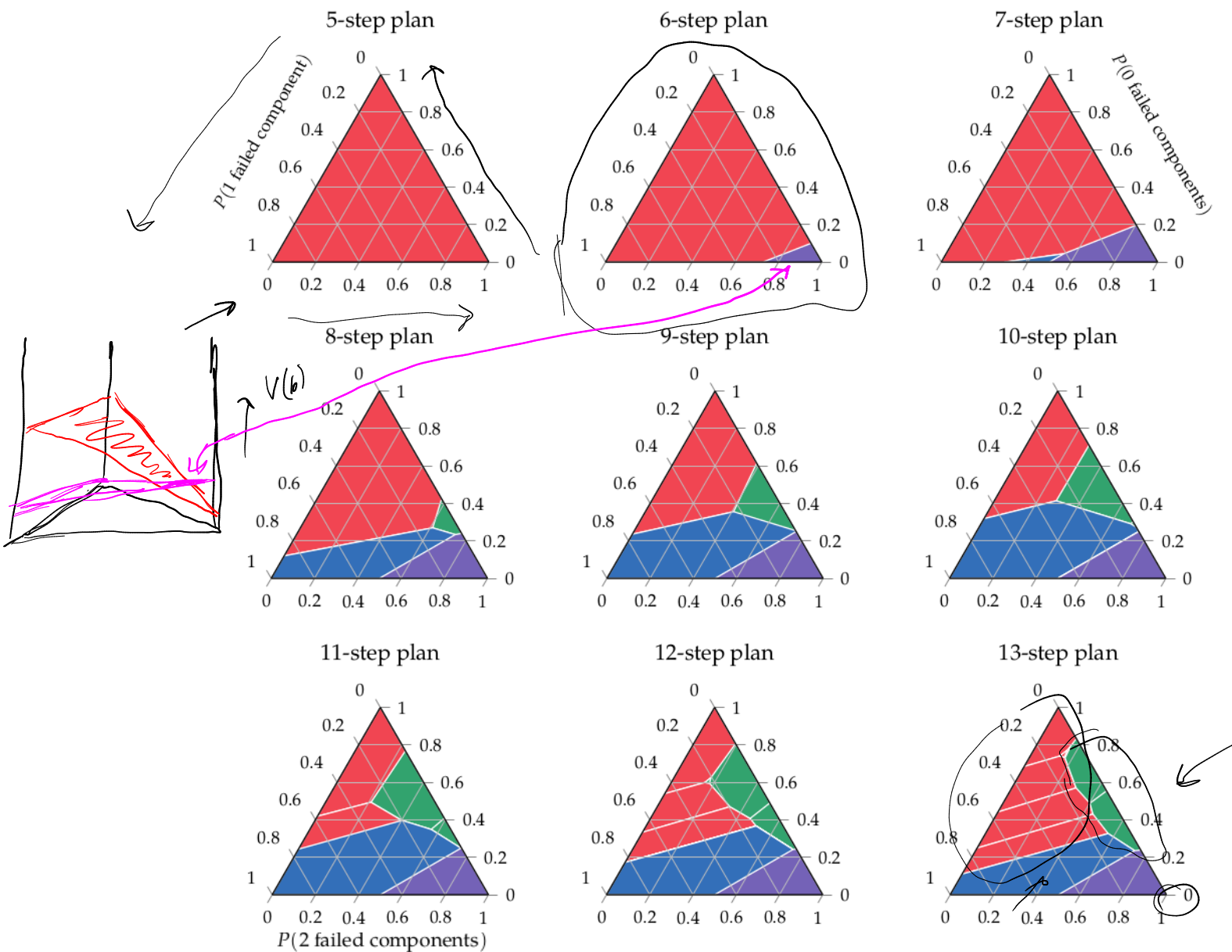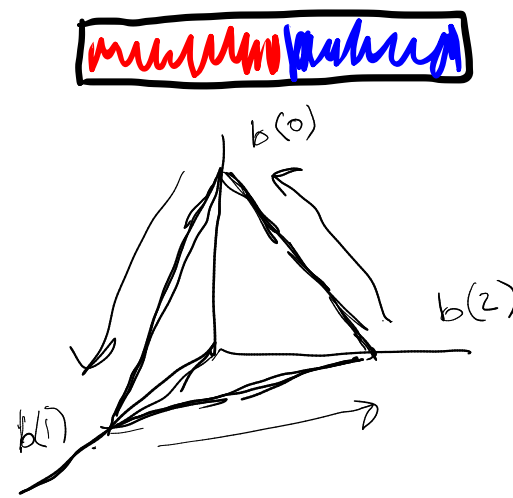
# Finite Horizon POMDP Value Iteration

# Finite Horizon POMDP Value Iteration



4

# Infinite-Horizon POMDP Lower Bound Improvement

$$\frac{R}{1-\gamma}$$

# Infinite-Horizon POMDP Lower Bound Improvement

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

$\Gamma \leftarrow$ blind lower bound

*always execute same action*

loop

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

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

A survey of point based POMDP solvers

backup

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

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

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

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

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

$$\Gamma^1 \oplus \Gamma^2 = \left\{ \alpha_1 + \alpha_2 : \alpha_1 \in \Gamma^1, \alpha_2 \in \Gamma^2 \right\}$$

# Point-Based Value Iteration (PBVI)
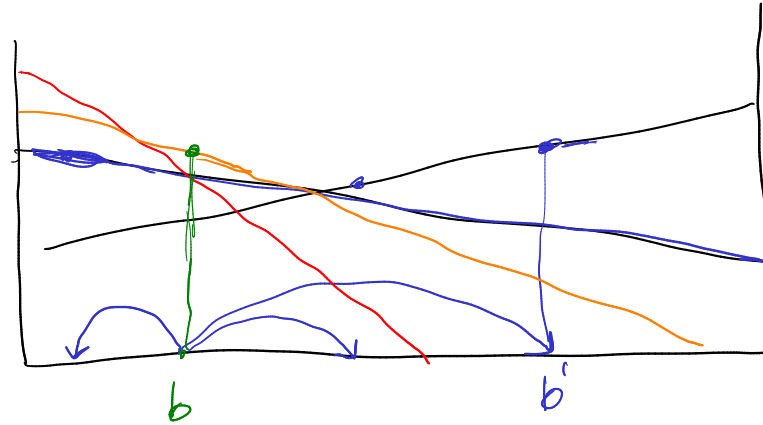
$point-$backup$(\Gamma, b)$

    for $a \in A$

        for $o \in O$

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

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

        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 $\underset{\alpha_a}{\operatorname{argmax}} \, \alpha_a^\top b$



$\mathcal{B}$

If we perform a backup for each $b \in \mathcal{B}$

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

6

# Original PBVI

how do we choose B

$B \leftarrow b_0$

loop

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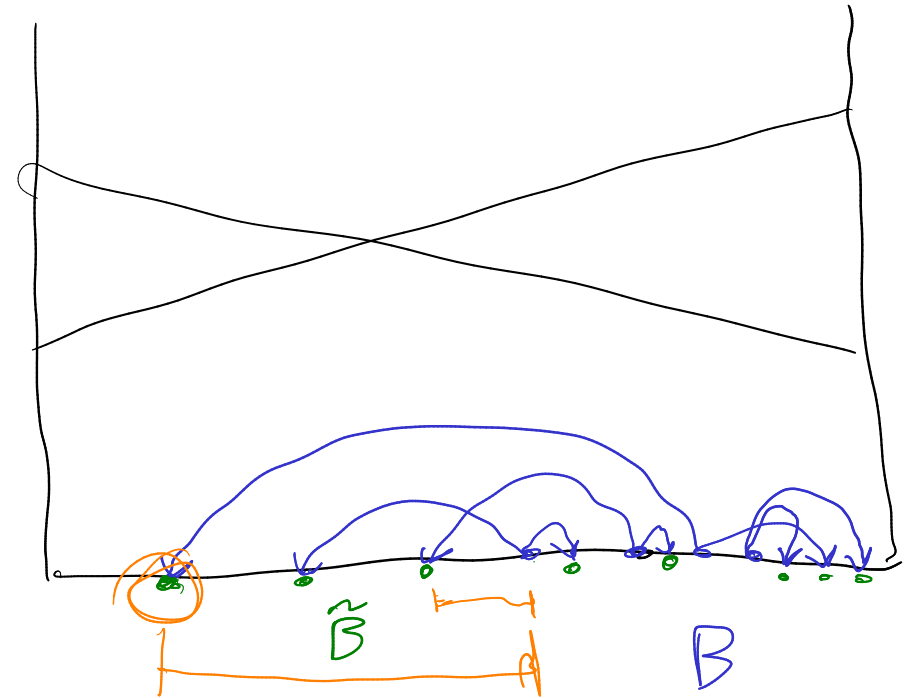
      $\Gamma \leftarrow \Gamma \cup \{\text{point\_backup}(\Gamma, b)\}$

   for $b \in B$

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

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

$B \leftarrow B \cup B'$

# PERSEUS: Randomly Selected Beliefs

Two Phases:

    1. Random Exploration
    2. Value Backup

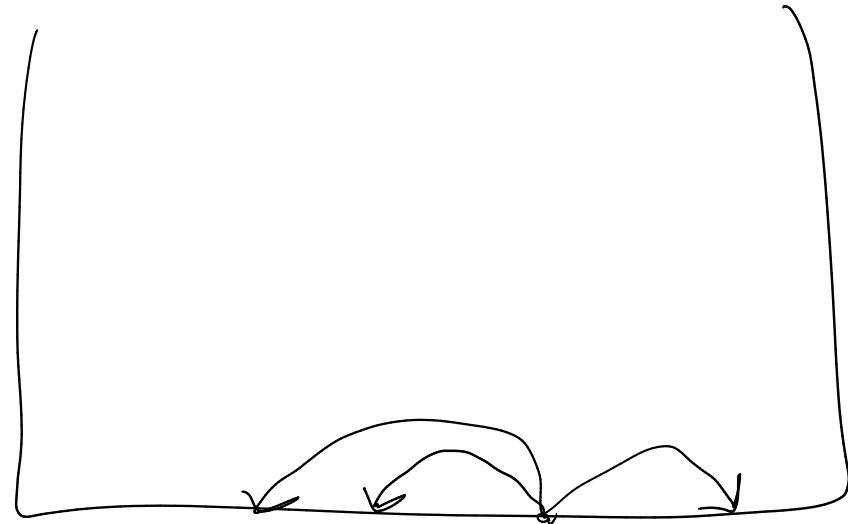Random Exploration:

$B \leftarrow \emptyset$

$b \leftarrow b_0$

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)

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

  explore$(b_0, 0)$

$\overline{V}(b)$      $\underline{V}(b)$

upper bound    lower bound
    for $b \in \mathcal{B}$

function explore(b, t)

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

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

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

  $WEU$

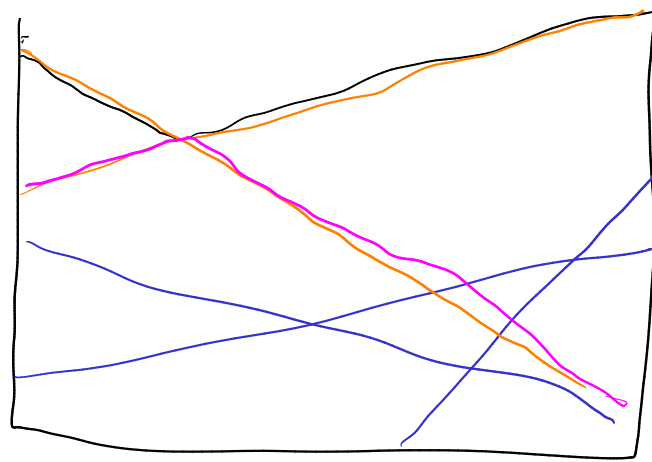  explore$(\tau(b, a^*, o^*), t+1)$

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

  $\overline{V}(b) = \underline{B}_b\left[\overline{V}(b)\right]$
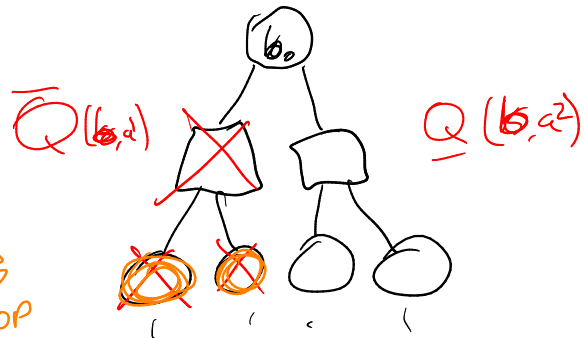
9

# Sawtooth Upper Bounds

# SARSOP

Successive Approximation of Reachable Space under Optimal Policies



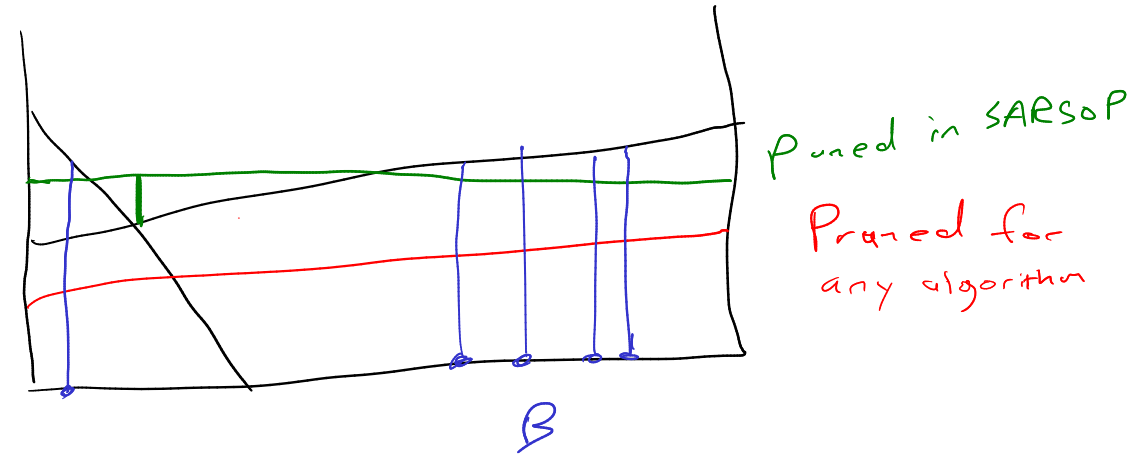Similar to HSVI

HSVI
$B \subset R$
↑
reachable

SARSOP
$B \subset R^*$
↑
reachable under
optimal policy

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

Not in B
for SARSOP

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

Pruned in SARSOP

Pruned for
any algorithm

B

Witness (α-vector value iteration) : ~20 states

SARSOP : 10,000 - 100,000 states

# Offline POMDP Algorithms

# Policy Graphs

# Monte Carlo Value Iteration (MCVI)