



# Offline POMDP Algorithms

Name		S	A	O
Survey of Point Based POMDP Solvers	Exact VI	D (10)	D	D
	PBVI	⋮	⋮	⋮
	Perseus	⋮	⋮	⋮
	HSVI	⋮	⋮	⋮
	SARSOP ← Use	D (1000-10,000)	⋮	⋮
MCVI		C	D	D/C

# Infinite-Horizon POMDP Value Iteration

constant -action  
policy

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

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

$\Gamma \leftarrow$  blind lower bound  
loop

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

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

backup

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

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

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

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

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

# Point-Based Value Iteration (PBVI)

```

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  $s \in S$ 

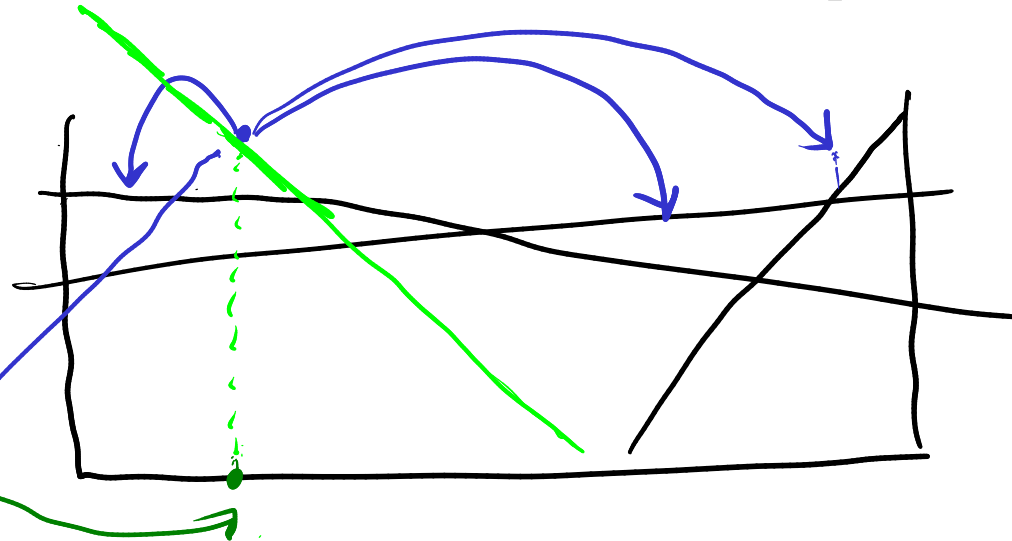
```

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

```

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

```



If we perform a point based backup for every  $b \in B$

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

How do we choose  $B$

# Original PBVI

$B \leftarrow b_0$

loop

for  $b \in B$

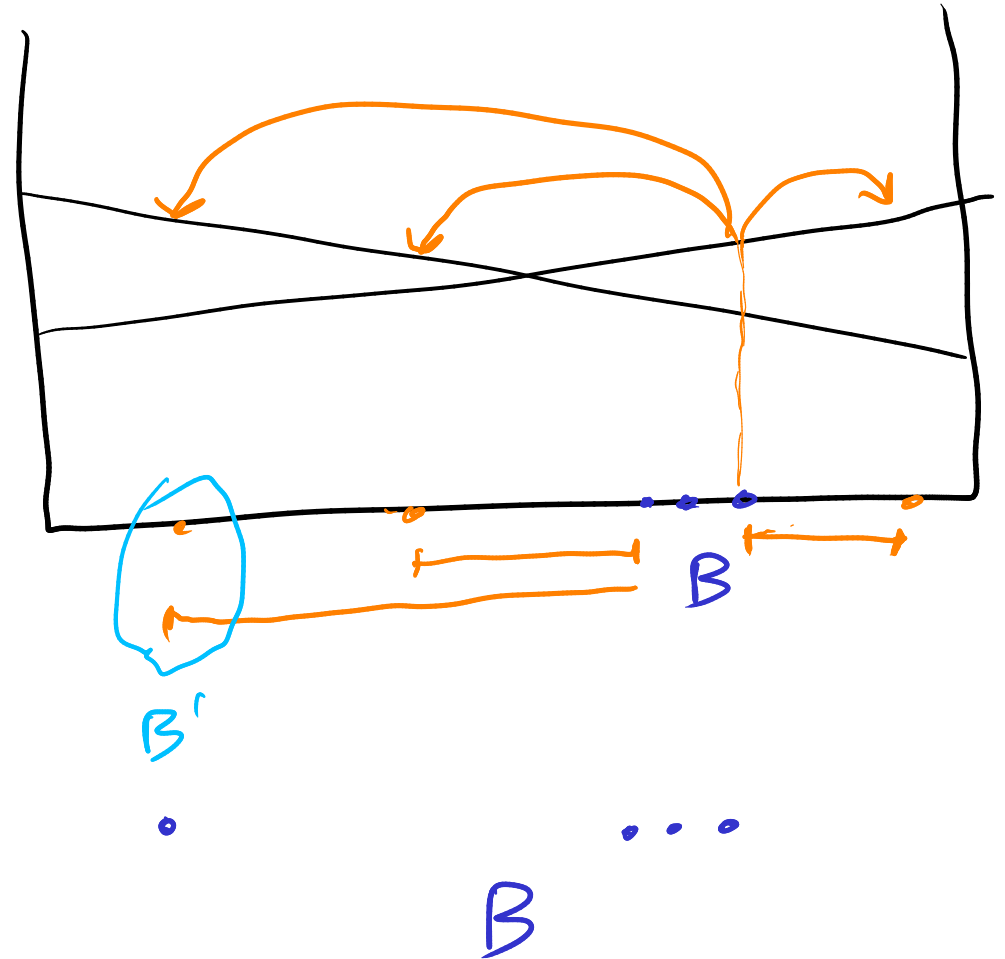
$\Gamma \leftarrow \Gamma \cup \{\text{point\_backup}(\Gamma, b)\}$

for  $b \in B$

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


$B' \leftarrow B' \cup \left\{ \underset{b' \in \tilde{B}}{\operatorname{argmax}} \underbrace{\|B, b'\|}_{\text{blue underline}} \right\}$

$B \leftarrow B \cup B'$



# PERSEUS: Randomly Selected Beliefs

Two Phases:

1. Random Exploration
2. Value Backup 

Random Exploration:

$$B \leftarrow \emptyset$$

$$b \leftarrow b_0$$

loop until  $|B| = n$

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

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

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

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

# Heuristic Search Value Iteration (HSVI)

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

  explore( $b_0, 0$ )

function explore( $b, t$ )

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

$a^* = \operatorname{argmax}_a \bar{Q}(b, a)$  ← maximize upper bound

$o^* = \operatorname{argmax}_o P(o \mid b, a) (\bar{V}(\tau(b, a^*, o)) - \underline{V}(\tau(b, a^*, o)) - \epsilon \gamma^t)$  ← Excess uncertainty

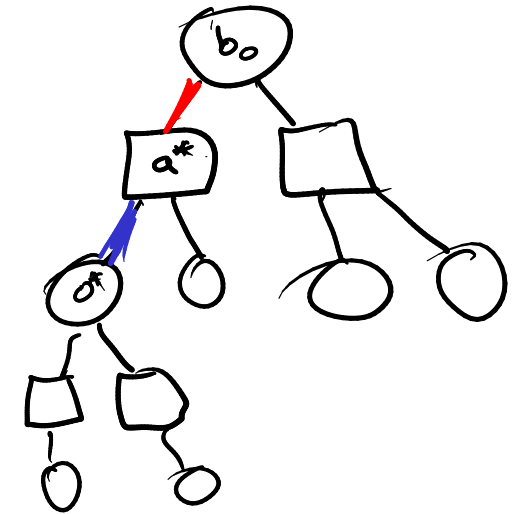
  explore( $\tau(b, a^*, o^*), t + 1$ )

$\underline{\Gamma} \leftarrow \underline{\Gamma} \cup \text{point\_backup}(\underline{\Gamma}, b)$  ← Lower Bound

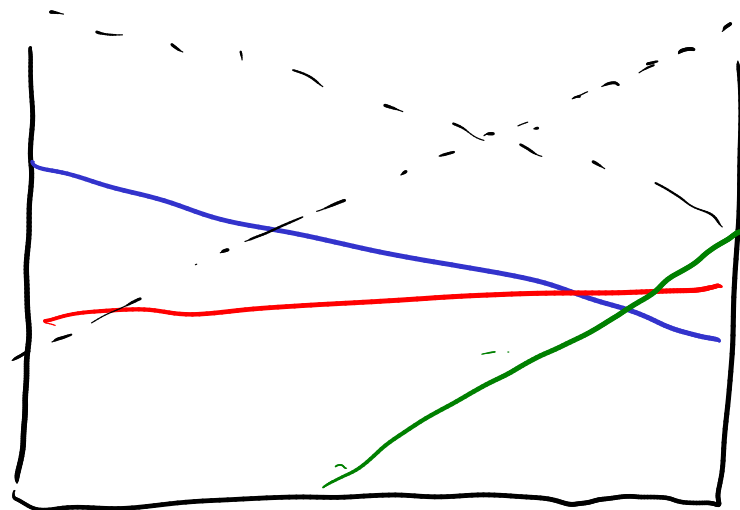
$\bar{V}(b) = B_b [\bar{V}(b)]$  ← Upper Bound

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

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



# Sawtooth Upper Bounds

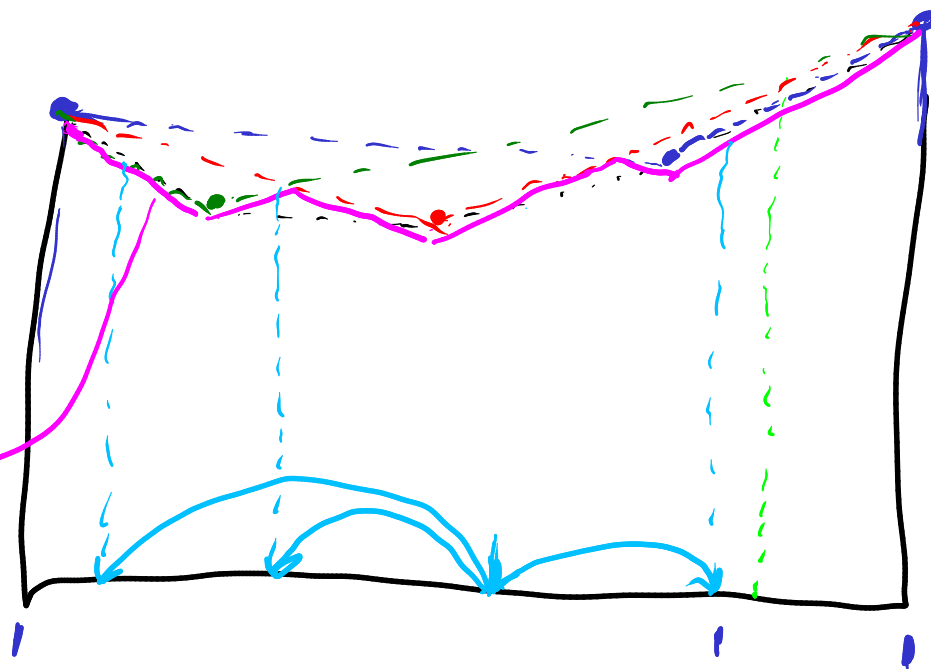


$$\min_{\alpha} \alpha^T b$$

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

$$B_b[\bar{V}(b)] = \max_a R(s, a)$$

$$+ \gamma \sum_o P(o|b, a) \bar{V}(\tau(b, a, o))$$



Finding the points to interpolate between requires solving a linear program



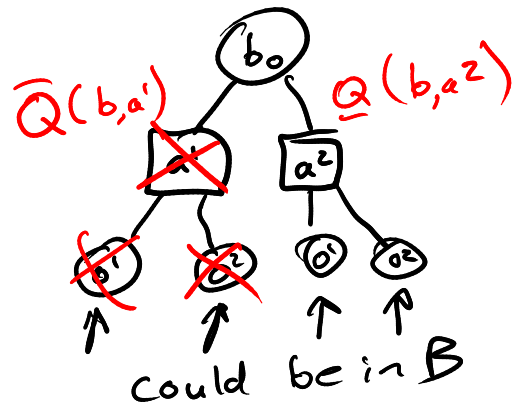
# SARSOP

Successive Approximation of Reachable Space under Optimal Policies

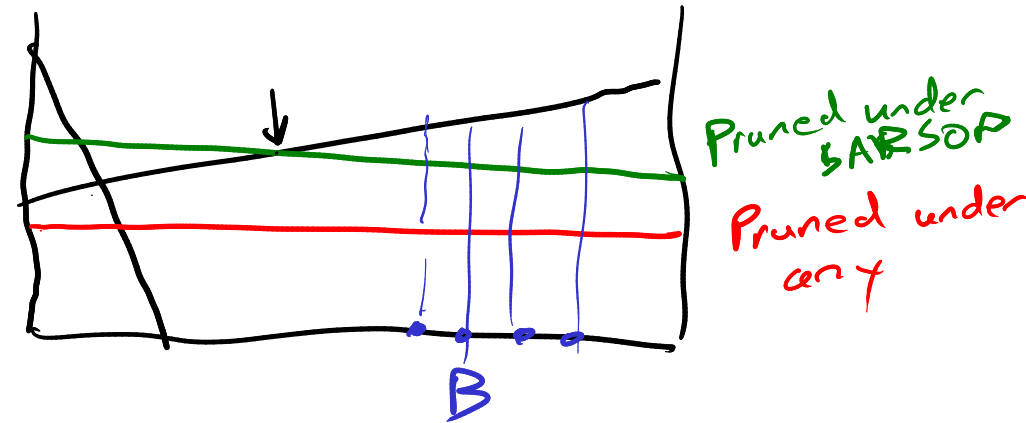
Similar to HSVI

HSVI  
 $B \subset R$   
 ↑  
 reachable

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



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



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

# Policy Graphs

# Monte Carlo Value Iteration (MCVI)