

# SARSOP

Hanna Kurniawati, David Hsu, and Wee Sun Lee

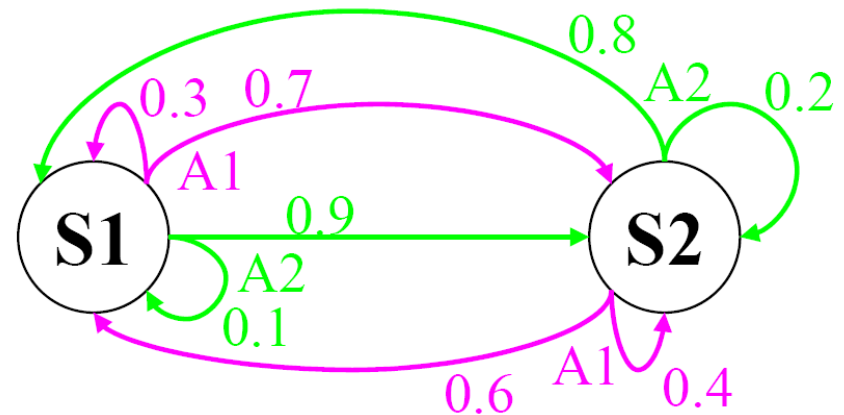
Presented by Geoff Hollinger

CS599

Spring, 2011

# From last week

- POMDPS
  - States, observations, actions, partial observability, Markov property
- Solution methods
  - Value iteration
  - Policy iteration
  - Point-based methods



# Key idea

- Can solve POMDPs by sampling the belief space

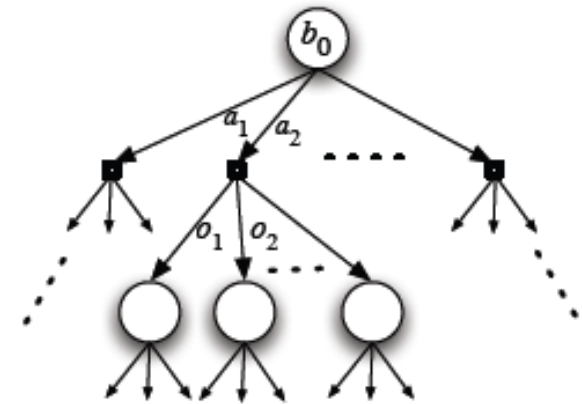


Fig. 2. The belief tree  $T_{\mathcal{R}}$  rooted at  $b_0$ .

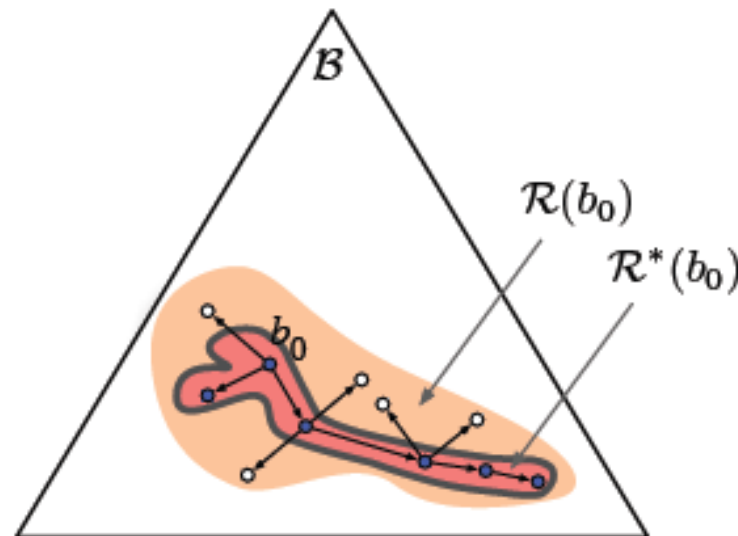


Fig. 1. Belief space  $\mathcal{B}$ , reachable space  $\mathcal{R}(b_0)$ , and optimally reachable space  $\mathcal{R}^*(b_0)$ . Note that  $\mathcal{R}^*(b_0) \subseteq \mathcal{R}(b_0) \subseteq \mathcal{B}$ .

# Covering number

- Covering number  $C(\delta)$  of a set  $S$  is the minimum number of balls of radius  $\delta$  needed to cover  $S$
- Measure of the “volume” of  $S$ .
- D. Hsu, “What makes some POMDPs problems easy to approximate,” NIPS, 2007.



# Covering number

- Key theorem from covering number
- SARSOP:
  - Want good performance when reachable belief space has a high covering number but the “optimally reachable belief space” has a low covering number

*Theorem 1:* For any  $b_0 \in \mathcal{B}$ , let  $\mathcal{C}(\delta)$  be the  $\delta$ -covering number of  $\mathcal{R}(b_0)$ . Given any constant  $\epsilon > 0$ , an approximation  $V(b_0)$  of  $V^*(b_0)$ , with error  $|V^*(b_0) - V(b_0)| \leq \epsilon$ , can be found in time

$$O \left( c \left( \frac{(1 - \gamma)^2 \epsilon}{4\gamma R_{\max}} \right)^2 \log_{\gamma} \frac{(1 - \gamma)\epsilon}{2R_{\max}} \right).$$

# Covering number

- Hardness of POMDPs

*Theorem 2:* Let  $b_0$  be any point in  $\mathcal{B}$  and  $\pi^*$  be an optimal policy. Given a constant  $\epsilon > 0$ , computing an approximation  $V(b_0)$  of  $V^*(b_0)$ , with error  $|V(b_0) - V^*(b_0)| \leq \epsilon |V^*(b_0)|$ , is NP-hard, even if the covering number of  $\mathcal{R}_{\pi^*}(b_0)$  is polynomial-sized.

- Not as hard when given a “cover”

*Theorem 3:* For any  $b_0 \in \mathcal{B}$  and any optimal policy  $\pi^*$ , given a proper  $\delta$ -cover  $C$  of  $\mathcal{R}_{\pi^*}(b_0)$  with  $\delta = \frac{(1-\gamma)^2 \epsilon}{2\gamma R_{\max}}$ , an approximation  $V(b_0)$  of  $V^*(b_0)$ , with error  $|V^*(b_0) - V(b_0)| \leq \epsilon$ , can be found in time

$$O\left(|C|^2 + |C| \log_{\gamma} \frac{(1-\gamma)\epsilon}{2R_{\max}}\right),$$

where  $|C|$  is the size of  $C$  and  $R_{\max} = \max_{s,a} |R(s,a)|$  is the maximum one-step reward.

# SARSOP

---

**Algorithm 1** SARSOP.

---

- 1: Initialize the set  $\Gamma$  of  $\alpha$ -vectors, representing the lower bound  $\underline{V}$  on the optimal value function  $V^*$ . Initialize the upper bound  $\overline{V}$  on  $V^*$ .
  - 2: Insert the initial belief point  $b_0$  as the root of the tree  $T_{\mathcal{R}}$ .
  - 3: **repeat**
  - 4:     SAMPLE( $T_{\mathcal{R}}, \Gamma$ ).
  - 5:     Choose a subset of nodes from  $T_{\mathcal{R}}$ . For each chosen node  $b$ , BACKUP( $T_{\mathcal{R}}, \Gamma, b$ ).
  - 6:     PRUNE( $T_{\mathcal{R}}, \Gamma$ ).
  - 7: **until** termination conditions are satisfied.
  - 8: **return**  $\Gamma$ .
-

# Backup operator

---

**Algorithm 2** Perform  $\alpha$ -vector backup at a node  $b$  of  $T_{\mathcal{R}}$ .

---

BACKUP( $T_{\mathcal{R}}, \Gamma, b$ )

- 1: For all  $a \in \mathcal{A}, o \in \mathcal{O}, \alpha_{a,o} \leftarrow \operatorname{argmax}_{\alpha \in \Gamma} (\alpha \cdot \tau(b, a, o))$ .
  - 2: For all  $a \in \mathcal{A}, s \in \mathcal{S}$ ,  
 $\alpha_a(s) \leftarrow R(s, a) + \gamma \sum_{o, s'} T(s, a, s') Z(s', a, o) \alpha_{a,o}(s')$ .
  - 3:  $\alpha' \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} (\alpha_a \cdot b)$
  - 4: Insert  $\alpha'$  into  $\Gamma$ .
- 

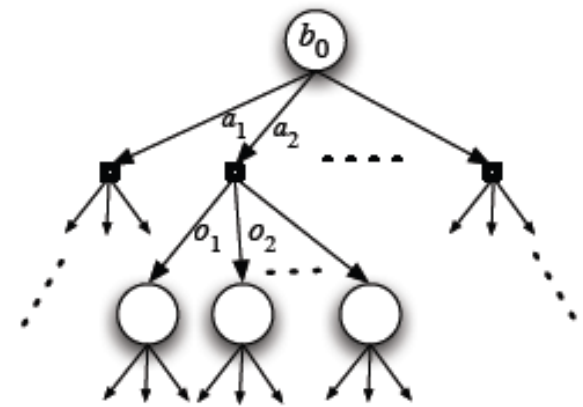
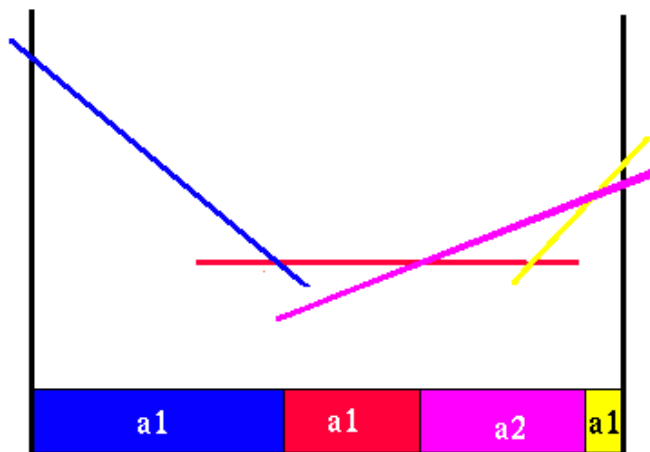


Fig. 2. The belief tree  $T_{\mathcal{R}}$  rooted at  $b_0$ .



# Sampling near $R^*$

- Want to sample the optimally reachable space
  - Maintain upper and lower bounds on cost function
    - Upper bound using Sawtooth bound (Hauskrecht, 2000)
    - Lower bound is the backup vectors
  - Choose action that has highest upper bound
  - Choose observation that minimizes gap in bounds
- Continue down the tree
  - Give preference to lower bound beyond gap of  $\gamma^{-t}\epsilon$
  - Predict  $V^*(b)$  at new node  $b$  and propagate towards root
    - Use belief clustering according to features (initial upper bound and entropy)
    - Pass lower bound down the tree and only propagate up when reached
  - Stop if gap termination reached (next slide)

# Gap termination

- Terminate and move up the tree by:
  - Passing an upper bound target level  $U$  down tree
  - Terminate if upper bound lower than  $U$
  - Allows the termination criterion to count all nodes (rather than just leaves)
- Ok, the upshot:
  - An effective sampling strategies that goes deep when needed, and avoids excessive sampling of belief space not in  $R^*$

# Sampling in SARSOP

- Selective deep sampling
  - Uses lower bound  $L$
- Termination condition
  - Uses upper bound  $U$

---

## Algorithm 3 Sampling near $\mathcal{R}^*$ .

---

SAMPLE( $T_{\mathcal{R}}, \Gamma$ )

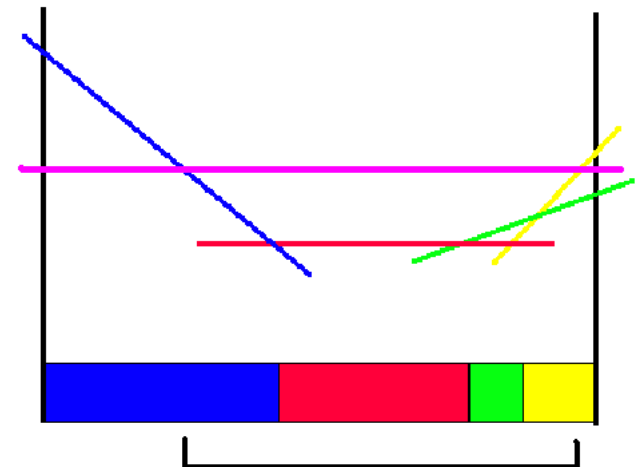
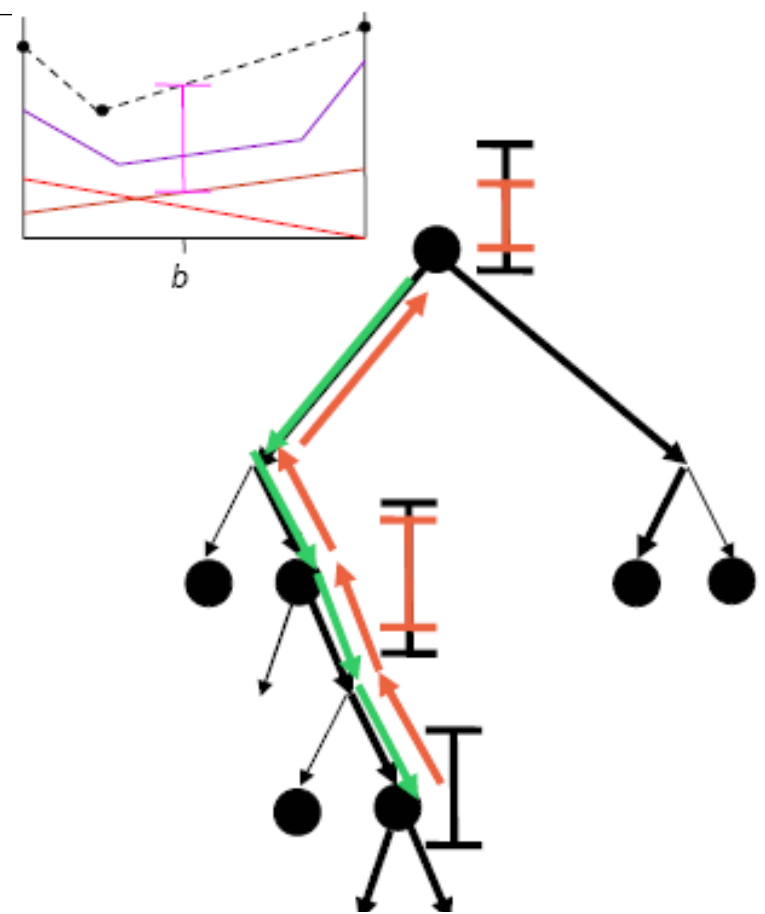
- 1: Set  $L$  to the current lower bound on the value function at the root  $b_0$  of  $T_{\mathcal{R}}$ . Set  $U$  to  $L + \epsilon$ , where  $\epsilon$  is the current target gap size at  $b_0$ .
- 2: SAMPLEPOINTS( $T_{\mathcal{R}}, \Gamma, b_0, L, U, \epsilon, 1$ ).

SAMPLEPOINTS( $T_{\mathcal{R}}, \Gamma, b, L, U, \epsilon, t$ ).

- 3: Let  $\hat{V}$  be the predicted value of  $V^*(b)$ .
  - 4: **if**  $\hat{V} \leq L$  and  $\bar{V}(b) \leq \max\{U, \underline{V}(b) + \epsilon\gamma^{-t}\}$  **then**
  - 5:     **return**
  - 6: **else**
  - 7:      $\underline{Q} \leftarrow \max_a \underline{Q}(b, a)$ .
  - 8:      $L' \leftarrow \max\{L, \underline{Q}\}$ .
  - 9:      $U' \leftarrow \max\{U, \underline{Q} + \gamma^{-t}\epsilon\}$ .
  - 10:      $a' \leftarrow \arg \max_a \bar{Q}(b, a)$ .
  - 11:      $o' \leftarrow \arg \max_o p(o|b, a') (\bar{V}(\tau(b, a', o)) - \underline{V}(\tau(b, a', o)))$ .
  - 12:     Calculate  $L_t$  so that  $L' = \sum_s R(s, a')b(s) + \gamma (p(o'|b, a')L_t + \sum_{o \neq o'} p(o|b, a')\underline{V}(\tau(b, a', o)))$ .
  - 13:     Calculate  $U_t$  so that  $U' = \sum_s R(s, a')b(s) + \gamma (p(o'|b, a')U_t + \sum_{o \neq o'} p(o|b, a')\bar{V}(\tau(b, a', o)))$ .
  - 14:      $b' \leftarrow \tau(b, a', o')$ .
  - 15:     Insert  $b'$  into  $T_{\mathcal{R}}$  as a child of  $b$ .
  - 16:     SAMPLEPOINTS( $T_{\mathcal{R}}, \Gamma, b', L_t, U_t, \epsilon, t + 1$ ).
-

# Pruning

- Node pruning
  - Prune if lower bound greater than upper bound
  - Prunes nodes that do not lie in  $R^*$  (not useful for optimal policy)
- Alpha vector pruning
  - Only prune if greater than  $\delta$  at ALL points along vector
  - Computed efficiently by intersection of hyperplanes



# SARSOP (again)

---

**Algorithm 1** SARSOP.

---

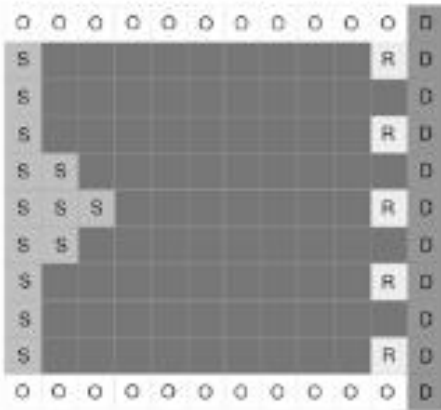
- 1: Initialize the set  $\Gamma$  of  $\alpha$ -vectors, representing the lower bound  $\underline{V}$  on the optimal value function  $V^*$ . Initialize the upper bound  $\overline{V}$  on  $V^*$ .
  - 2: Insert the initial belief point  $b_0$  as the root of the tree  $T_{\mathcal{R}}$ .
  - 3: **repeat**
  - 4:     SAMPLE( $T_{\mathcal{R}}, \Gamma$ ).
  - 5:     Choose a subset of nodes from  $T_{\mathcal{R}}$ . For each chosen node  $b$ , BACKUP( $T_{\mathcal{R}}, \Gamma, b$ ).
  - 6:     PRUNE( $T_{\mathcal{R}}, \Gamma$ ).
  - 7: **until** termination conditions are satisfied.
  - 8: **return**  $\Gamma$ .
-

# APPL software

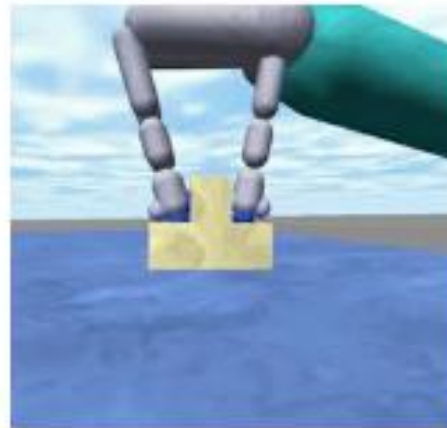
- **POMDP software**
  - **Implements SARSOP**
  - **Same format as HSVI2 software**
- **Relatively easy to get to work (I've used it)**

<http://motion.comp.nus.edu.sg/projects/pomdp/pomdp.html>.

# Applications



navigation



grasping



target tracking



integrated exploration

# Applications

TABLE I  
PERFORMANCE COMPARISON.

	Reward	Time (s)
<b>Underwater Navigation,</b> $ S =2,653,  A =6,  O =103$		
SARSOP	$722.59 \pm 1.30$	72
HSV12	$721.45 \pm 0.75$	720
<b>Grasping</b> $ S =1,253,  A =6,  O =96$		
SARSOP	$320.00 \pm 0.16$	8
HSV12	$319.88 \pm 0.14$	60
<b>Integrated Exploration</b> $ S =15,517,  A =8,  O =1,015$		
SARSOP	$(1.58 \pm 0.03) \times 10^6$	5,400
HSV12	$(1.41 \pm 0.02) \times 10^6$	5,400
	$(1.43 \pm 0.02) \times 10^6$	7,200
<b>Rock Sample (7,8)</b> $ S =12,545,  A =13,  O =2$		
SARSOP	$21.27 \pm 0.13$	400
HSV12	$21.27 \pm 0.09$	250
<b>Tag</b> $ S =870,  A =5,  O =30$		
SARSOP	$-6.13 \pm 0.12$	6
HSV12	$-7.43 \pm 0.11$	6
	$-6.40 \pm 0.10$	7,200
<b>Homecare</b> $ S =5,408,  A =9,  O =928$		
SARSOP	$16.86 \pm 0.45$	960
HSV12	$16.88 \pm 0.37$	2,880



# Tag problem

- Need to catch a non-adversarial target
- Motion model is Markovian
- Reward gained when moving to the same state as the target
- Target's state is “partially observed”

