

# L17 – Solving POMDPs Offline

AIMA 17.5, KAlg 20.5-6; 21.4-5

## What you should know after this lecture

- How to specify a full POMDP policy
- What an  $\alpha$ -vector is
- How to use value iteration to compute a value function
- Point-based value iteration methods find policy concentrated on belief space reachable under optimal policy
- New online belief-space planning methods can be smarter than expectimax and work in complicated domains

## Recall POMDP definitions

MDP with added observation process

- MDP has
  - a set of states  $\mathcal{S}$
  - a set of actions  $\mathcal{A}$
  - transition model
$$T(s, a, s') = P(S_{t+1} = s' \mid S_t = s, A_t = a)$$
  - reward function  $R : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$
  - discount factor  $\gamma$  (and possibly horizon  $T$ )
- POMDP adds
  - a set of observations  $\mathcal{O}$
  - observation model
$$O(s', a, o) = P(O_{t+1} = o \mid A_t = a, S_{t+1} = s')$$
- optionally

## Recall Bayes filter belief update

Given previous belief  $\mathbf{b}$ , action  $\mathbf{a}$  and observation  $\mathbf{o}$ , what is the new belief  $\mathbf{b}' = \text{bf}(\mathbf{b}, \mathbf{a}, \mathbf{o})$ ?

$$\begin{aligned}\text{bf}(\mathbf{b}, \mathbf{a}, \mathbf{o})(s') &= P(S_{t+1} = s' \mid A_t = \mathbf{a}, O_{t+1} = \mathbf{o}, B_t = \mathbf{b}) \\ &= \frac{1}{\eta} P(O_{t+1} = \mathbf{o} \mid S_{t+1} = s', A_t = \mathbf{a}, B_t = \mathbf{b}) \\ &\quad P(S_{t+1} = s' \mid B_t = \mathbf{b}, A_t = \mathbf{a}) \\ &= \frac{1}{\eta} O(s', \mathbf{a}, \mathbf{o}) \sum_s P(S_{t+1} = s' \mid B_t = \mathbf{b}, A_t = \mathbf{a}, S_t = s) \\ &\quad P(S_t = s \mid B_t = \mathbf{b}, A_t = \mathbf{a}) \\ &= \frac{1}{\eta} O(s', \mathbf{a}, \mathbf{o}) \sum_s T(s, \mathbf{a}, s') \mathbf{b}(s)\end{aligned}$$

where  $\eta = P(\mathbf{o} \mid \mathbf{a}, \mathbf{b}) = \sum_{\tilde{s}} O(\tilde{s}, \mathbf{a}, \mathbf{o}) \sum_s T(s, \mathbf{a}, \tilde{s}) \mathbf{b}(s)$

## Value of a policy tree: $\alpha$ vector

Let's start by computing the value of executing a policy tree  $\pi$  in a known starting state  $s$ :  $V^\pi(s)$ .

- Base case when  $\pi$  is single node ( $H = 1$ ):

$$V^\pi(s) = R(s, \pi.a)$$

- Recursive case: depending on what state  $s'$  we transition to and what observation we get (which depends on the state), we will execute one of our subtrees ( $\pi(o)$ ) in  $s'$ :

$$V^\pi(s) = R(s, \pi.a) + \gamma \sum_{s'} T(s, \pi.a, s') \sum_o O(s', \pi.a, o) V^{\pi(o)}(s')$$

Let

$$\alpha^\pi = [V^\pi(s^1), \dots, V^\pi(s^{|\mathcal{S}|})]$$

Then value at a belief is

$$V^\pi(\mathbf{b}) = \sum_s \mathbf{b}(s) V^\pi(s) = \mathbf{b} \cdot \alpha^\pi$$

## Infinite-horizon discounted case

- Optimal value function is convex
- It can be piecewise linear or curved! Curve can arise in the limit of infinitely many pieces.
- Value iteration algorithm still works, iteratively computing sets of  $\alpha$  vectors.
- Terminate when the change in the maximum difference between subsequent value functions becomes small.
- Cool (advanced topic): if the optimal value function has finitely many pieces, then there is a finite-state machine controller that is optimal!

## POMDP “backup”

What is the value, at belief state  $\mathbf{b}$ , of taking action  $\mathbf{a}$  and then, for each  $\mathbf{o} \in \mathcal{O}$ , if we get observation  $\mathbf{o}$ , continuing with a policy whose value is  $\alpha_{\mathbf{o}}$ ?

$$\begin{aligned}V_{\mathbf{a}}(\mathbf{b}) &= \left( \sum_s \mathbf{b}(s)R(s, \mathbf{a}) \right) + \gamma \sum_{\mathbf{o}} P(\mathbf{o} | \mathbf{b}, \mathbf{a}) \alpha_{\mathbf{o}} \cdot \mathbf{b}f(\mathbf{b}, \mathbf{a}, \mathbf{o}) \\&= \left( \sum_s \mathbf{b}(s)R(s, \mathbf{a}) \right) + \gamma \sum_{\mathbf{o}} P(\mathbf{o} | \mathbf{b}, \mathbf{a}) \sum_{s'} (\alpha_{\mathbf{o}}(s') \mathbf{b}f(\mathbf{b}, \mathbf{a}, \mathbf{o})(s')) \\&= \left( \sum_s \mathbf{b}(s)R(s, \mathbf{a}) \right) + \gamma \sum_{\mathbf{o}} P(\mathbf{o} | \mathbf{b}, \mathbf{a}) \sum_{s'} \alpha_{\mathbf{o}}(s') \frac{O(s', \mathbf{a}, \mathbf{o}) \sum_s T(s, \mathbf{a}, s') \mathbf{b}(s)}{P(\mathbf{o} | \mathbf{b}, \mathbf{a})} \\&= \sum_s \mathbf{b}(s) \left( R(s, \mathbf{a}) + \gamma \sum_{\mathbf{o}} \sum_{s'} \alpha_{\mathbf{o}}(s') O(s', \mathbf{a}, \mathbf{o}) T(s, \mathbf{a}, s') \right) \\&= \mathbf{b} \cdot \alpha'\end{aligned}$$

We get a new alpha vector!  $\text{BACKUP}(\mathbf{a}, [\alpha_1, \dots, \alpha_{|\mathcal{O}|}]) = \alpha'$

# POMDP Value Iteration

POMDP-VI( $\mathcal{S}, \mathcal{A}, \mathbb{T}, \mathbb{R}, \mathcal{O}, \mathbf{O}, \gamma$ )

```
1   $\Gamma_0 = \{\mathbb{R}(\cdot, a^1), \mathbb{R}(\cdot, a^2), \dots, \mathbb{R}(\cdot, a^{|\mathcal{A}|})\}$            // H=1  $\alpha$  vectors
2   $t = 1; \Delta = \infty$ 
3  while  $\Delta > \epsilon$            // Stop when max change in value funs is  $< \epsilon$ 
4       $\Gamma_t = \{\}$ 
5      for  $a \in \mathcal{A}$            // Try all combinations of subtrees
6          for  $(\alpha_1, \dots, \alpha_{|\mathcal{O}|}) \in \text{COMBINATIONS}(\Gamma_{t-1}, |\mathcal{O}|)$ 
7               $\alpha = \text{BACKUP}(a, [\alpha_1, \dots, \alpha_{|\mathcal{O}|}])$ 
8               $\Gamma_t = \Gamma_t \cup \alpha$ 
9          // Ideally, prune out dominated  $\alpha$  from  $\Gamma_t$ 
10      $\Delta = \max_b ((\max_{\alpha_t \in \Gamma_t} \alpha_t \cdot b) - (\max_{\alpha_{t-1} \in \Gamma_{t-1}} \alpha_{t-1} \cdot b))$ 
11      $t = t + 1$ 
12 return  $\Gamma_t$ 
```

## Variations on value iteration

- Guaranteed to converge to optimum but can be very slow because there may be many tiny little “facets” to the value function
- Idea: sample specific points in belief space to control where we spend our computational / approximation effort.

## Point-based value iteration

Point-based backup computes new  $\alpha$  vector that is guaranteed to be an improvement at  $b$  (and probably elsewhere).

PB-BACKUP( $POMDP, \Gamma, b$ )

- 1 For  $a \in \mathcal{A}$
- 2     For  $o \in \mathcal{O}$
- 3          $\alpha_{a,o} = \operatorname{argmax}_{\alpha \in \Gamma} \alpha \cdot \mathbf{bf}(b, a, o)$
- 4      $\alpha_a = \text{BACKUP}(POMDP, a, [\alpha_{a,1}, \dots, \alpha_{a,|\mathcal{O}|}])$
- 5 **return**  $\operatorname{argmax}_a \alpha_a \cdot b$

Randomly sampling  $b$  will converge to optimal  $V$  but slow.

PBVI( $POMDP$ )

- 1  $\Gamma = \{R(\cdot, a^1), R(\cdot, a^2), \dots, R(\cdot, a^{|\mathcal{A}|})\}$
- 2 **while** not tired
- 3      $b = \text{SAMPLE-B}(\Gamma)$
- 4      $\Gamma = \Gamma \cup \{\text{PB-BACKUP}(POMDP, \Gamma, b)\}$
- 5 **return**  $\Gamma$

# SARSOP

Idea: Use  $b_0$  and only try to estimate  $V$  well on belief states reachable via optimal policy.  $V_\Gamma$  is a lower bound on  $V^*$ .

SARSOP( $POMDP, b_0, \epsilon$ )

- 1  $\Gamma = \{R(\cdot, a^1), R(\cdot, a^2), \dots, R(\cdot, a^{|\mathcal{A}|})\}$
- 2 Initialize upper bound  $V_{up}$  (e.g. with  $Q_{mdp}$ )
- 3 Initialize partial expectimax tree  $T$  with root  $b_0$
- 4 **while**  $|V_\Gamma(b_0) - V_{up}(b_0)| > \epsilon$
- 5      $b_1, \dots, b_k = \text{SAMPLE-PATH}(\Gamma, T)$
- 6     For  $b_i \in b_1, \dots, b_k$
- 7          $\Gamma = \Gamma \cup \{\text{PB-BACKUP}(\Gamma, b_i)\}$
- 8         Update  $V_{up}$  at  $b_i$
- 9         Update  $T(b_i)$
- 10     PRUNE( $\Gamma$ )
- 11 **return**  $\Gamma$

Hanna Kurniawati, David Hsu, Wee Sun Lee, SARSOP: Efficient point-based

POMDP planning by approximating optimally reachable belief spaces.

## SARSOP sampling: intuition only

To sample a path, start at root of  $T$  and generate a path. At a node  $b$

- Select  $a$  that maximizes upper bound  $Q_{\text{up}}(b, a)$
- Select  $o$  that maximizes  $|V_{\text{up}}(\mathbf{bf}(b, a, o)) - V_{\Gamma}(\mathbf{bf}(b, a, o))|$

We try to stay in high-value parts of the state space and to sample in places where our bounds are far apart.

Terminate a sample trajectory if the difference between upper and lower bounds is such that it will have little effect on the gap at  $b_0$ .

## Modern online solution methods

We looked at an approximate online solution method: most likely observation. Plan under assumption of MLO, replan when we get a different observation.

MLO can be bad when there's a possible very bad outcome of an action that is not highly likely. It will not be "revealed" by the MLO and we will ignore the downside risk.

Can also do expectimax or sparse sampling on the belief MDP. MCTS offers more focused search:

- POMCP algorithm (Silver et al): labels nodes with  $o$ , a histories rather than beliefs — allows approximate belief representations such as particle sets
- DESPOT algorithm (Ye et al): Uses cleverer sampling (some ideas from SARSOP) and variance reduction techniques to be more efficient than POMCP

# Next time

- Real RL is a POMDP!
- We'll start with Bandit problems