

## 20 *Exact Belief State Planning*

The objective in a POMDP is to choose actions that maximize the accumulation of reward while interacting with the environment. In contrast with MDPs, states are not directly observable, requiring the agent to use its past history of actions and observations to inform a belief. As discussed in the previous chapter, beliefs can be represented as probability distributions over states. There are different approaches for computing an optimal policy that maps beliefs to actions given models of the transitions, observations, and rewards.<sup>1</sup> One approach is to convert a POMDP into an MDP and apply dynamic programming. Other approaches include representing policies as conditional plans or as piecewise linear value functions over the belief space. The chapter concludes with an algorithm for computing an optimal policy that is analogous to value iteration for MDPs.

<sup>1</sup> A discussion of exact solution methods is provided by L. P. Kaelbling, M. L. Littman, and A. R. Cassandra, "Planning and Acting in Partially Observable Stochastic Domains," *Artificial Intelligence*, vol. 101, no. 1–2, pp. 99–134, 1998.

### 20.1 *Belief-State Markov Decision Processes*

Any POMDP can be viewed as an MDP that uses beliefs as states, also called a *belief-state MDP*.<sup>2</sup> The state space of a belief-state MDP is the set of all beliefs  $\mathcal{B}$ . The action space is identical to that of the POMDP.

The reward function for a belief-state MDP depends on the belief and action taken. It is simply the expected value of the reward. For a discrete state-space, it is given by

$$R(b, a) = \sum_s R(s, a) b(s) \quad (20.1)$$

<sup>2</sup> K. J. Åström, "Optimal Control of Markov Processes with Incomplete State Information," *Journal of Mathematical Analysis and Applications*, vol. 10, no. 1, pp. 174–205, 1965.

If the state and observation spaces are discrete, the belief-state transition function for a belief-state MDP is given by

$$T(b' | b, a) = P(b' | b, a) \quad (20.2)$$

$$= \sum_o P(b' | b, a, o) P(o | b, a) \quad (20.3)$$

$$= \sum_o P(b' | b, a, o) \sum_s P(o | b, a, s) P(s | b, a) \quad (20.4)$$

$$= \sum_o P(b' | b, a, o) \sum_s P(o | b, a, s) b(s) \quad (20.5)$$

$$= \sum_o P(b' | b, a, o) \sum_{s'} \sum_s P(o | b, a, s, s') P(s' | b, s, a) b(s) \quad (20.6)$$

$$= \sum_o (b' = \text{Update}(b, a, o)) \sum_{s'} O(o | a, s') \sum_s T(s' | s, a) b(s) \quad (20.7)$$

In equation (20.7),  $\text{Update}(b, a, o)$  returns the updated belief using the deterministic process discussed in the previous chapter.<sup>3</sup> For continuous problems, we replace the summations with integrals.

Solving belief-state MDPs is challenging because the state space is continuous. We can use the approximate dynamic programming techniques presented in earlier chapters, but we can often do better by taking advantage of the structure of the belief-state MDP, as will be discussed in the remainder of this chapter.

## 20.2 Conditional Plans

There are a number of ways to represent policies for POMDPs. One approach is to use a *conditional plan* represented as a tree. Figure 20.1 shows an example of a three-step conditional plan with binary action and observation spaces. The nodes correspond to belief states. The edges are annotated with observations, and the nodes are annotated with actions. If we have a plan  $\pi$ , the action associated with the root is denoted as  $\pi()$  and the subplan associated with observation  $o$  is denoted as  $\pi(o)$ . Algorithm 20.1 provides an implementation of this.

A conditional plan tells us what to do in response to our observations up to the horizon represented by the tree. To execute a conditional plan, we start with the root node and execute the action associated with it. We proceed down the tree according to our observations, taking the actions associated with the nodes through which we pass.

<sup>3</sup> As a reminder, we use the convention where a logical statement in parentheses is treated numerically as 1 when true and 0 when false.

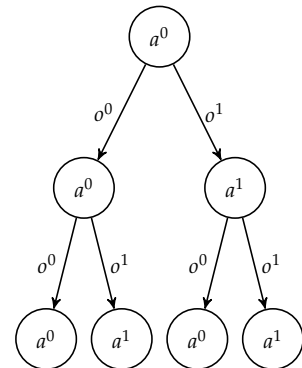


Figure 20.1. A three-step conditional plan.

```

struct ConditionalPlan
  a      # action to take at root
  subplans # dictionary mapping observations to subplans
end

ConditionalPlan(a) = ConditionalPlan(a, Dict())

( $\pi$ ::ConditionalPlan)() =  $\pi$ .a
( $\pi$ ::ConditionalPlan)(o) =  $\pi$ .subplans[o]

```

Algorithm 20.1. The conditional plan data structure consisting of an action and a mapping from observations to subplans. The `subplans` field is a `Dict` from observations to conditional plans. For convenience, we have created a special constructor for plans that consist of a single node.

Suppose we have a conditional plan  $\pi$ , and we want to compute its expected utility when starting from state  $s$ . This computation can be done recursively:

$$U^\pi(s) = R(s, \pi()) + \gamma \left[ \sum_{s'} T(s' | s, \pi()) \sum_o O(o | \pi(), s') U^{\pi(o)}(s') \right] \quad (20.8)$$

An implementation for this procedure is given in algorithm 20.2.

```

function lookahead( $\mathcal{P}$ ::POMDP, U, s, a)
  S, O, T, O, R,  $\gamma$  =  $\mathcal{P}$ .S,  $\mathcal{P}$ .O,  $\mathcal{P}$ .T,  $\mathcal{P}$ .O,  $\mathcal{P}$ .R,  $\mathcal{P}$ . $\gamma$ 
  u' = sum(T(s,a,s')*sum(O(a,s',o)*U(o,s') for o in O) for s' in S)
  return R(s,a) +  $\gamma$ *u'
end

function evaluate_plan( $\mathcal{P}$ ::POMDP,  $\pi$ ::ConditionalPlan, s)
  U(o,s') = evaluate_plan( $\mathcal{P}$ ,  $\pi$ (o), s')
  return isempty( $\pi$ .subplans) ?  $\mathcal{P}$ .R(s, $\pi$ ()) : lookahead( $\mathcal{P}$ , U, s,  $\pi$ ())
end

```

Algorithm 20.2. A method for evaluating a conditional plan  $\pi$  for MDP  $\mathcal{P}$  starting at state  $s$ . Plans are represented as tuples consisting of an action and a dictionary mapping observations to subplans.

We can compute the utility of our belief  $b$  as follows:

$$U^\pi(b) = \sum_s b(s) U^\pi(s) \quad (20.9)$$

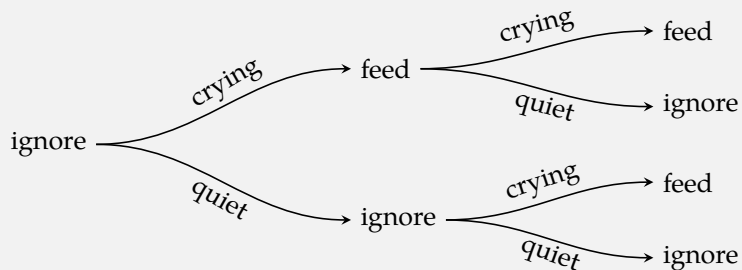
Example 20.1 shows how to compute the utility associated with a three-step conditional plan.

Now that we have a way to evaluate conditional plans up to a horizon  $h$ , we can compute the optimal  $h$ -step value function:

$$U^*(b) = \max_{\pi} U^\pi(b) \quad (20.10)$$

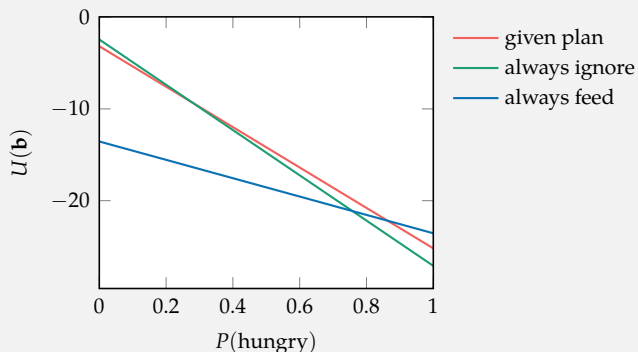
An optimal action can be generated from the action associated with the root of a maximizing  $\pi$ .

Consider the following three-step conditional plan for the crying baby problem:



In this plan, we begin by ignoring the baby. If we observe any crying, we feed the baby. If we do not observe any crying, we ignore the baby. Our third action again feeds if there is crying.

The expected utility for this plan in belief space is plotted alongside a three-step plan that always feeds the baby and one that always ignores the baby.



We find that the given plan is not universally better than either always ignoring or always feeding the baby.

Example 20.1. A conditional plan for the three-step crying baby problem (appendix F.7), evaluated and compared to two simpler conditional plans.

Solving an  $h$ -step POMDP by directly enumerating all  $h$ -step conditional plans is generally computationally intractable, as shown in figure 20.2. There are  $(|\mathcal{O}|^h - 1)/(|\mathcal{O}| - 1)$  nodes in an  $h$ -step plan. In general, any action can be inserted into any node, resulting in  $|\mathcal{A}|^{(|\mathcal{O}|^h - 1)/(|\mathcal{O}| - 1)}$  possible  $h$ -step plans. This exponential growth means that enumerating over all plans is intractable even for modest values of  $h$ . As will be discussed later in this chapter, there are alternatives to explicitly enumerating over all possible plans.

### 20.3 Alpha Vectors

We can rewrite equation (20.9) in vector form:

$$U^\pi(\mathbf{b}) = \sum_s b(s)U^\pi(s) = \boldsymbol{\alpha}_\pi^\top \mathbf{b} \quad (20.11)$$

The vector  $\boldsymbol{\alpha}_\pi$ , called an *alpha vector*, contains the expected utility under plan  $\pi$  for each state. As with belief vectors, alpha vectors have dimension  $|\mathcal{S}|$ . Unlike beliefs, the components in alpha vectors represent utilities, not probability masses. Algorithm 20.3 shows how to compute an alpha vector.

```
function alphavector( $\mathcal{P}$ ::POMDP,  $\pi$ ::ConditionalPlan)
    return [evaluate_plan( $\mathcal{P}$ ,  $\pi$ ,  $s$ ) for  $s$  in  $\mathcal{P}.\mathcal{S}$ ]
end
```

Each alpha vector defines a hyperplane in belief space. The optimal value function given in equation (20.11) is the maximum over these hyperplanes:

$$U^*(\mathbf{b}) = \max_{\pi} \boldsymbol{\alpha}_\pi^\top \mathbf{b} \quad (20.12)$$

making the value function piecewise-linear and convex.<sup>4</sup>

An alternative to using a conditional plan to represent a policy is to use a set of alpha vectors  $\Gamma$ , each annotated with an action. Although it is not practical, one way to generate set  $\Gamma$  is to enumerate the set of  $h$ -step conditional plans and then compute their alpha vectors. The action associated with an alpha vector is the action at the root of the associated conditional plan. We execute a policy represented by  $\Gamma$  by updating our belief state and performing the action associated with the dominating alpha vector at the new belief  $\mathbf{b}$ . The dominating alpha vector  $\boldsymbol{\alpha}$  at  $\mathbf{b}$  is the one that maximizes  $\boldsymbol{\alpha}^\top \mathbf{b}$ . This strategy can be used to select actions

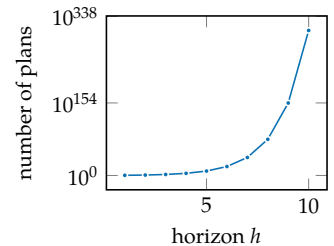


Figure 20.2. Even for small POMDPs with only two actions and two observations, the number of possible plans grows extremely quickly with the planning horizon. We can often significantly prune the set of alpha vectors at each iteration and only consider far fewer plans.

Algorithm 20.3. We can generate an alpha vector from a conditional plan by calling `evaluate_plan` from all possible initial states.

<sup>4</sup>The optimal value function for continuous-state POMDPs is also convex, as can be seen by approximating the POMDP through state space discretization and taking the limit as the number of discrete states approaches infinity.

beyond the horizon of the original conditional plans. Algorithm 20.4 provides an implementation.

```

struct AlphaVectorPolicy
   $\mathcal{P}$  # POMDP problem
   $\Gamma$  # alpha vectors
   $\mathbf{a}$  # actions associated with alpha vectors
end

function utility( $\pi$ ::AlphaVectorPolicy,  $\mathbf{b}$ )
  return maximum( $\alpha \cdot \mathbf{b}$  for  $\alpha$  in  $\pi.\Gamma$ )
end

function ( $\pi$ ::AlphaVectorPolicy)( $\mathbf{b}$ )
   $i$  = argmax( $[\alpha \cdot \mathbf{b}$  for  $\alpha$  in  $\pi.\Gamma$ ])
  return  $\pi.\mathbf{a}[i]$ 
end

```

Algorithm 20.4. An alpha vector policy is defined in terms of a set of alpha vectors  $\Gamma$  and an array of associated actions  $\mathbf{a}$ . Given the current belief  $\mathbf{b}$ , it will find the alpha vector that gives the highest value at that belief point. It will return the associated action.

If we use *one-step lookahead*, we do not have to keep track of the actions associated with the alpha vectors in  $\Gamma$ . The one-step lookahead action from belief  $\mathbf{b}$  using the value function represented by  $\Gamma$ , denoted as  $U^\Gamma$ , is

$$\pi^\Gamma(\mathbf{b}) = \arg \max_a \left[ R(\mathbf{b}, a) + \gamma \sum_o P(o | \mathbf{b}, a) U^\Gamma(\text{Update}(\mathbf{b}, a, o)) \right] \quad (20.13)$$

where

$$P(o | \mathbf{b}, a) = \sum_s P(o | s, a) b(s) \quad (20.14)$$

$$P(o | s, a) = \sum_{s'} T(s' | s, a) O(o | s', a) \quad (20.15)$$

Algorithm 20.5 provides an implementation of this. Example 20.2 demonstrates using one-step lookahead on the crying baby problem.

## 20.4 Pruning

If we have a collection of alpha vectors  $\Gamma$ , we may want to *prune* alpha vectors that do not contribute to our representation of the value function or plans that are not optimal for any belief. Removing such alpha vectors or plans can improve computational efficiency. We can check whether an alpha vector  $\alpha$  is *dominated* by

```

function lookahead( $\mathcal{P}$ ::POMDP, U, b::Vector, a)
     $S, \mathcal{O}, T, \mathcal{O}, R, \gamma = \mathcal{P}.S, \mathcal{P}.\mathcal{O}, \mathcal{P}.T, \mathcal{P}.\mathcal{O}, \mathcal{P}.R, \mathcal{P}.\gamma$ 
     $r = \text{sum}(R(s,a)*b[i] \text{ for } (i,s) \text{ in enumerate}(S))$ 
     $\text{Posa}(o,s,a) = \text{sum}(\mathcal{O}(a,s',o)*T(s,a,s') \text{ for } s' \text{ in } S)$ 
     $\text{Poba}(o,b,a) = \text{sum}(b[i]*\text{Posa}(o,s,a) \text{ for } (i,s) \text{ in enumerate}(S))$ 
    return  $r + \gamma*\text{sum}(\text{Poba}(o,b,a)*U(\text{update}(b, \mathcal{P}, a, o)) \text{ for } o \text{ in } \mathcal{O})$ 
end

function greedy( $\mathcal{P}$ ::POMDP, U, b::Vector)
    u, a = findmax(a→lookahead( $\mathcal{P}$ , U, b, a),  $\mathcal{P}.A$ )
    return (a=a, u=u)
end

struct LookaheadAlphaVectorPolicy
     $\mathcal{P}$  # POMDP problem
     $\Gamma$  # alpha vectors
end

function utility( $\pi$ ::LookaheadAlphaVectorPolicy, b)
    return maximum( $\alpha \cdot b$  for  $\alpha$  in  $\pi.\Gamma$ )
end

function greedy( $\pi$ , b)
    U(b) = utility( $\pi$ , b)
    return greedy( $\pi.\mathcal{P}$ , U, b)
end

( $\pi$ ::LookaheadAlphaVectorPolicy)(b) = greedy( $\pi$ , b).a

```

Algorithm 20.5. A policy represented by a set of alpha vectors  $\Gamma$ . It uses one-step lookahead to produce an optimal action and associated utility. Equation (20.13) is used to compute the lookahead.

Consider using one-step lookahead on the crying baby problem with a value function given by the alpha vectors  $[-3.7, -15]$  and  $[-2, -21]$ . Suppose that our current belief is  $b = [0.5, 0.5]$ , meaning that we believe it is equally likely the baby is hungry as not hungry. We apply equation (20.13)

$$\begin{aligned}
 & R(b, \text{feed}) = -10 \\
 & \left. \begin{aligned}
 & \gamma P(\text{crying} \mid b, \text{feed}) U(\text{Update}(b, \text{feed}, \text{crying})) = -0.18 \\
 & \gamma P(\text{quiet} \mid b, \text{feed}) U(\text{Update}(b, \text{feed}, \text{quiet})) = -1.62 \\
 & \rightarrow Q(b, \text{feed}) = -11.8
 \end{aligned} \right\} \\
 & R(b, \text{ignore}) = -5 \\
 & \left. \begin{aligned}
 & \gamma P(\text{crying} \mid b, \text{ignore}) U(\text{Update}(b, \text{ignore}, \text{crying})) = -6.09 \\
 & \gamma P(\text{quiet} \mid b, \text{ignore}) U(\text{Update}(b, \text{ignore}, \text{quiet})) = -2.81 \\
 & \rightarrow Q(b, \text{ignore}) = -13.9
 \end{aligned} \right\} \\
 & R(b, \text{sing}) = -5.5 \\
 & \left. \begin{aligned}
 & \gamma P(\text{crying} \mid b, \text{sing}) U(\text{Update}(b, \text{sing}, \text{crying})) = -6.68 \\
 & \gamma P(\text{quiet} \mid b, \text{sing}) U(\text{Update}(b, \text{sing}, \text{quiet})) = -1.85 \\
 & \rightarrow Q(b, \text{sing}) = -14.0
 \end{aligned} \right\}
 \end{aligned}$$

We use  $Q(b, a)$  to represent the action value function from a belief state. The policy predicts that feeding the baby will result in the highest expected utility, so it takes that action.

Example 20.2. Applying a lookahead policy to the crying baby problem.



the alpha vectors in a set  $\Gamma$  by solving a linear program to maximize the utility gap  $\delta$  that vector achieves over all other vectors:<sup>5</sup>

$$\begin{aligned}
 & \underset{\delta, \mathbf{b}}{\text{maximize}} && \delta \\
 & \text{subject to} && \mathbf{b} \geq \mathbf{0} \\
 & && \mathbf{1}^\top \mathbf{b} = 1 \\
 & && \boldsymbol{\alpha}^\top \mathbf{b} \geq \boldsymbol{\alpha}'^\top \mathbf{b} + \delta, \quad \boldsymbol{\alpha}' \in \Gamma
 \end{aligned} \tag{20.16}$$

<sup>5</sup> Constraints of the form  $\mathbf{a} \geq \mathbf{b}$  are elementwise. That is, we mean  $a_i \geq b_i$  for all  $i$ .

The first two constraints ensure that  $\mathbf{b}$  is a categorical distribution, and the final set of constraints ensures that we find a belief vector for which  $\boldsymbol{\alpha}$  has a higher expected reward than all alpha vectors in  $\Gamma$ . If, after solving the linear program, the utility gap  $\delta$  is negative, then  $\boldsymbol{\alpha}$  is dominated. If  $\delta$  is positive, then  $\boldsymbol{\alpha}$  is not dominated and  $\mathbf{b}$  is a belief at which  $\boldsymbol{\alpha}$  is not dominated. Algorithm 20.6 provides an implementation for solving equation (20.16) to determine a belief, if one exists, where  $\delta$  is most positive.

```

function find_maximal_belief( $\alpha$ ,  $\Gamma$ )
    m = length( $\alpha$ )
    if isempty( $\Gamma$ )
        return fill(1/m, m) # arbitrary belief
    end
    model = Model(GLPK.Optimizer)
    @variable(model,  $\delta$ )
    @variable(model, b[i=1:m]  $\geq$  0)
    @constraint(model, sum(b) == 1.0)
    for a in  $\Gamma$ 
        @constraint(model, ( $\alpha$ -a)·b  $\geq$   $\delta$ )
    end
    @objective(model, Max,  $\delta$ )
    optimize!(model)
    return value( $\delta$ ) > 0 ? value.(b) : nothing
end

```

Algorithm 20.6. A method for finding the belief vector  $\mathbf{b}$  for which the alpha vector  $\boldsymbol{\alpha}$  improves the most compared to the set of alpha vectors  $\Gamma$ . Nothing is returned if no such belief exists. The packages JuMP.jl and GLPK.jl provide a mathematical optimization framework and a solver for linear programs, respectively.

Algorithm 20.7 shows a procedure that uses algorithm 20.6 to find the dominating alpha vectors in a set  $\Gamma$ . Initially, all the alpha vectors are candidates for being dominating. We then choose one of these candidates and determine the belief  $\mathbf{b}$  where the candidate leads to the greatest improvement in value compared to all other alpha vectors in the dominating set. If the candidate does not bring improvement, we remove it from the set. If it does bring improvement, we move an alpha vector from the candidate set that brings the greatest improvement

at  $b$  to the dominating set. The process continues until there are no longer any candidates. We can prune away any alpha vectors and associated conditional plans that are not dominating at any belief point. Example 20.3 demonstrates pruning on the crying baby problem.

```
function find_dominating( $\Gamma$ )
  n = length( $\Gamma$ )
  candidates, dominating = trues(n), falses(n)
  while any(candidates)
    i = findfirst(candidates)
    b = find_maximal_belief( $\Gamma$ [i],  $\Gamma$ [dominating])
    if b === nothing
      candidates[i] = false
    else
      k = argmax([candidates[j] ? b. $\Gamma$ [j] : -Inf for j in 1:n])
      candidates[k], dominating[k] = false, true
    end
  end
  return dominating
end

function prune(plans,  $\Gamma$ )
  d = find_dominating( $\Gamma$ )
  return (plans[d],  $\Gamma$ [d])
end
```

Algorithm 20.7. A method for pruning dominated alpha vectors and associated plans. The `find_dominating` function identifies all the dominating alpha vectors in set  $\Gamma$ . It uses binary vectors `candidates` and `dominating` to track which alpha vectors are candidates for inclusion in the dominating set and which are currently in the dominating set, respectively.

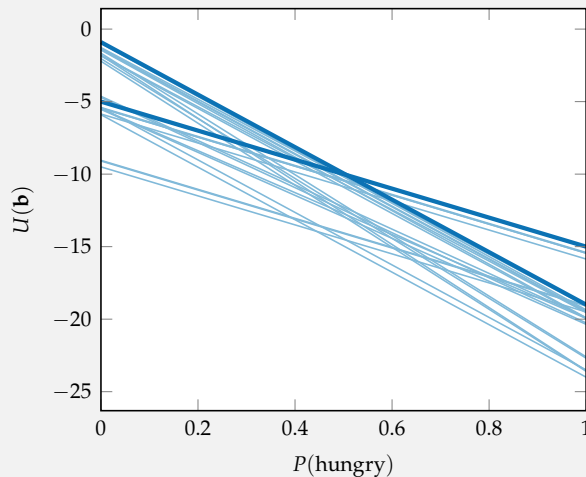
## 20.5 Value Iteration

The value iteration algorithm for MDPs can be adapted for POMDPs.<sup>6</sup> POMDP *value iteration* (algorithm 20.8) begins by constructing all one-step plans. We prune any plans that are never optimal for any initial belief. Then, we expand all combinations of one-step plans to produce two-step plans. Again, we prune any suboptimal plans from consideration. This procedure of alternating between expansion and pruning is repeated until the desired horizon is reached. Figure 20.3 demonstrates value iteration on the crying baby problem.

<sup>6</sup> This section describes a version of value iteration in terms of conditional plans and alpha vectors. For a version that only uses alpha vectors, see A. R. Cassandra, M. L. Littman, and N. L. Zhang, “Incremental Pruning: A Simple, Fast, Exact Method for Partially Observable Markov Decision Processes,” in *Conference on Uncertainty in Artificial Intelligence (UAI)*, 1997.

We can construct all two-step plans for the crying baby problem. There are  $3^3 = 27$  such plans.

The expected utility for each plan in belief space is plotted below. We find that two plans dominate all others. These dominating plans are the only ones that need to be considered as subplans for optimal three-step plans.



Example 20.3. The expected utility over the belief space for all two-step plans for the crying baby problem (appendix F.7). The thick lines are optimal for some beliefs, whereas the thin lines are dominated.

```
function value_iteration( $\mathcal{P}$ ::POMDP, k_max)
     $S, \mathcal{A}, R = \mathcal{P}.S, \mathcal{P}.A, \mathcal{P}.R$ 
    plans = [ConditionalPlan(a) for a in  $\mathcal{A}$ ]
     $\Gamma = [[R(s,a) \text{ for } s \text{ in } S] \text{ for } a \text{ in } \mathcal{A}]$ 
    plans,  $\Gamma = \text{prune}(\text{plans}, \Gamma)$ 
    for k in 2:k_max
        plans,  $\Gamma = \text{expand}(\text{plans}, \Gamma, \mathcal{P})$ 
        plans,  $\Gamma = \text{prune}(\text{plans}, \Gamma)$ 
    end
    return (plans,  $\Gamma$ )
end

function solve( $M$ ::ValueIteration,  $\mathcal{P}$ ::POMDP)
    plans,  $\Gamma = \text{value\_iteration}(\mathcal{P}, M.k\_max)$ 
    return LookaheadAlphaVectorPolicy( $\mathcal{P}, \Gamma$ )
end
```

Algorithm 20.8. Value iteration for POMDPs, which finds the dominating  $h$ -step plans for a finite horizon POMDP of horizon  $k\_max$  by iteratively constructing optimal subplans. The `ValueIteration` structure is the same as what was defined in algorithm 7.8 in the context of MDPs.

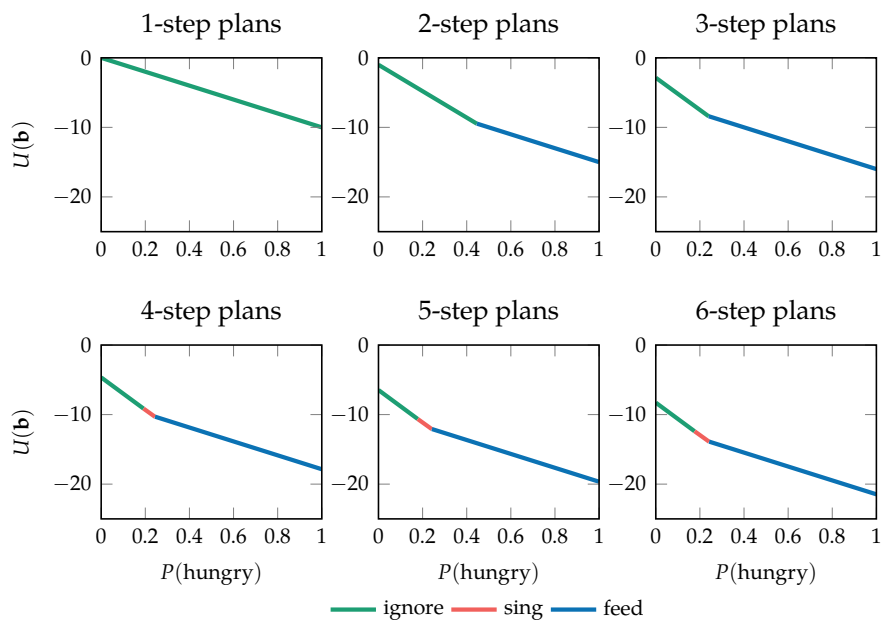


Figure 20.3. POMDP value iteration used to find the optimal value function for the crying baby problem to various horizons.

The expansion step (algorithm 20.9) in this process constructs all possible  $(k + 1)$ -step plans from a set of  $k$ -step plans. New plans can be constructed using a new first action and all possible combinations of the  $k$ -step plans as subplans, as shown in figure 20.4. While plans can also be extended by adding actions to the ends of subplans, top-level expansion allows alpha vectors constructed for the  $k$ -step plans to be used to efficiently construct alpha vectors for the  $(k + 1)$ -step plans.

Computing the alpha vector associated with a plan  $\pi$  from a set of alpha vectors associated with its subplans can be done as follows. We use  $\alpha_o$  to represent the alpha vector associated with subplan  $\pi(o)$ . The alpha vector associated with  $\pi$  is then

$$\alpha(s) = R(s, \pi()) + \gamma \sum_{s'} T(s' | s, \pi()) \sum_o O(o | \pi(), s') \alpha_o(s') \quad (20.17)$$

Even for relatively simple problems to shallow depths, computing alpha vectors from subplans in this way is much more efficient than computing them from scratch, as in algorithm 20.2.

## 20.6 Linear Policies

As discussed in section 19.3, the belief state in a problem with linear Gaussian dynamics can be represented by a Gaussian distribution,  $\mathcal{N}(\mu_b, \Sigma_b)$ . If the reward function is quadratic, then it can be shown that the optimal policy can be computed exactly offline using a process that is often called *linear quadratic Gaussian (LQG)* control. The optimal action is obtained in an identical manner as in section 7.8, but the  $\mu_b$  computed using the linear Gaussian filter is treated as the true state.<sup>7</sup> With each observation, we simply use the filter to update our  $\mu_b$  and obtain an optimal action by multiplying  $\mu_b$  with the policy matrix from algorithm 7.11. Example 20.4 demonstrates this process.

## 20.7 Summary

- Exact solutions for POMDPs typically can be obtained only for finite horizon discrete POMDPs.

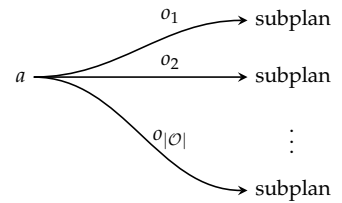


Figure 20.4. A  $(k + 1)$ -step plan can be constructed using a new initial action leading to any combination of  $k$ -step subplans.

<sup>7</sup> Our ability to simply use the mean of the distribution is another instance of the *certainty equivalence* principle, originally introduced in section 7.8.

```

function ConditionalPlan( $\mathcal{P}$ ::POMDP, a, plans)
    subplans = Dict( $o \Rightarrow \pi$  for ( $o, \pi$ ) in zip( $\mathcal{P}.O$ , plans))
    return ConditionalPlan(a, subplans)
end

function combine_lookahead( $\mathcal{P}$ ::POMDP, s, a,  $\Gamma_o$ )
     $S, O, T, O, R, \gamma = \mathcal{P}.S, \mathcal{P}.O, \mathcal{P}.T, \mathcal{P}.O, \mathcal{P}.R, \mathcal{P}.\gamma$ 
     $U'(s', i) = \text{sum}(O(a, s', o) * \alpha[i] \text{ for } (o, \alpha) \text{ in zip}(O, \Gamma_o))$ 
    return  $R(s, a) + \gamma * \text{sum}(T(s, a, s') * U'(s', i) \text{ for } (i, s') \text{ in enumerate}(S))$ 
end

function combine_alphavector( $\mathcal{P}$ ::POMDP, a,  $\Gamma_o$ )
    return [combine_lookahead( $\mathcal{P}$ , s, a,  $\Gamma_o$ ) for s in  $\mathcal{P}.S$ ]
end

function expand(plans,  $\Gamma, \mathcal{P}$ )
     $S, \mathcal{A}, O, T, O, R = \mathcal{P}.S, \mathcal{P}.\mathcal{A}, \mathcal{P}.O, \mathcal{P}.T, \mathcal{P}.O, \mathcal{P}.R$ 
    plans',  $\Gamma' = [], []$ 
    for a in  $\mathcal{A}$ 
        # iterate over all possible mappings from observations to plans
        for inds in product([eachindex(plans) for o in  $O$ ]...)
             $\pi_o = \text{plans}[[\text{inds}...]]$ 
             $\Gamma_o = \Gamma[[\text{inds}...]]$ 
             $\pi = \text{ConditionalPlan}(\mathcal{P}, a, \pi_o)$ 
             $\alpha = \text{combine_alphavector}(\mathcal{P}, a, \Gamma_o)$ 
            push!(plans',  $\pi$ )
            push!( $\Gamma'$ ,  $\alpha$ )
        end
    end
    return (plans',  $\Gamma'$ )
end

```

Algorithm 20.9. The expansion step in value iteration, which constructs all  $(k + 1)$ -step conditional plans and associated alpha vectors from a set of  $k$ -step conditional plans and alpha vectors. The way that we combine alpha vectors of subplans follows equation (20.17).

Consider a satellite navigating in two dimensions, neglecting gravity, drag, and other external forces. The satellite can use its thrusters to accelerate in any direction with linear dynamics:

$$\begin{bmatrix} x \\ y \\ \dot{x} \\ \dot{y} \end{bmatrix} \leftarrow \begin{bmatrix} 1 & 0 & \Delta t & 0 \\ 0 & 1 & 0 & \Delta t \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ \dot{x} \\ \dot{y} \end{bmatrix} + \begin{bmatrix} \frac{1}{2}\Delta t^2 & 0 \\ 0 & \frac{1}{2}\Delta t^2 \\ \Delta t & 0 \\ 0 & \Delta t \end{bmatrix} \begin{bmatrix} \ddot{x} \\ \ddot{y} \end{bmatrix} + \boldsymbol{\epsilon}$$

where  $\Delta t$  is the duration of a time step and  $\boldsymbol{\epsilon}$  is zero-mean Gaussian noise with covariance  $\Delta t/20\mathbf{I}$ .

We seek to place the satellite in its orbital slot at the origin, while minimizing fuel use. Our quadratic reward function is

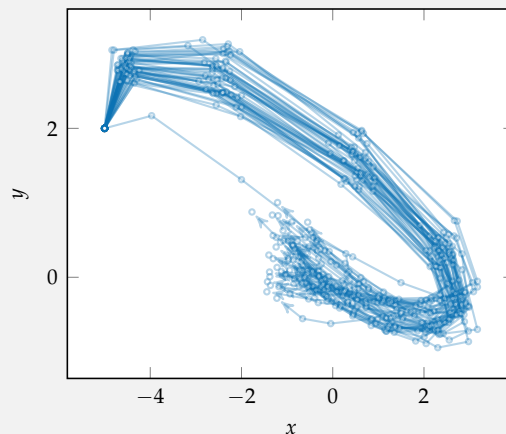
$$R(\mathbf{s}, \mathbf{a}) = -\mathbf{s}^\top \begin{bmatrix} \mathbf{I}_{2 \times 2} & \mathbf{0}_{2 \times 2} \\ \mathbf{0}_{2 \times 2} & \mathbf{0}_{2 \times 2} \end{bmatrix} \mathbf{s} - 2\mathbf{a}^\top \mathbf{a}$$

The satellite's sensors measure its position according to:

$$\mathbf{o} = \begin{bmatrix} \mathbf{I}_{2 \times 2} & \mathbf{0}_{2 \times 2} \end{bmatrix} \mathbf{s} + \boldsymbol{\varepsilon}$$

where  $\boldsymbol{\varepsilon}$  is zero-mean Gaussian noise with covariance  $\Delta t/10\mathbf{I}$ .

Here are 50 trajectories from 10-step rollouts using the optimal policy for  $\Delta t = 1$  and a Kalman filter to track the belief. In each case, the satellite was started at  $\mathbf{s} = \boldsymbol{\mu}_b = [-5, 2, 0, 1]$  with  $\boldsymbol{\Sigma}_b = [\mathbf{I}\mathbf{0}; \mathbf{0}\mathbf{0.25}\mathbf{I}]$ .



Example 20.4. An optimal policy used for a POMDP with linear Gaussian dynamics and quadratic reward.

- Policies for these problems can be represented as conditional plans, which are trees that describe the actions to take based on the observations.
- Alpha vectors contain the expected utility when starting from different states and following a particular conditional plan.
- Alpha vectors can also serve as an alternative representation of a POMDP policy.
- POMDP value iteration can avoid the computational burden of enumerating all conditional plans by iteratively computing subplans and pruning those that are suboptimal.
- Linear Gaussian problems with quadratic reward can be solved exactly using methods very similar to those derived for the fully observable case.

## 20.8 Exercises

**Exercise 20.1.** Can every POMDP be framed as an MDP?

*Solution:* Yes. Any POMDP can equivalently be viewed as a belief-state MDP whose state space is the space of beliefs in the POMDP, whose action space is the same as that of the POMDP and whose transition function is given by equation (20.2).

**Exercise 20.2.** What are the alpha vectors for the one-step crying baby problem (appendix F.7)? Are all the available actions dominant?

*Solution:* There are three one-step conditional plans, one for each action, resulting in three alpha vectors. The optimal one-step policy must choose between these actions, given the current belief. The one-step alpha vectors for a POMDP can be obtained from the optimal one-step belief value function:

$$U^*(b) = \max_a \sum_s b(s)R(s, a)$$

Feeding the baby yields an expected reward:

$$\begin{aligned} R(\text{hungry}, \text{feed})P(\text{hungry}) + R(\text{sated}, \text{feed})P(\text{sated}) \\ = -15P(\text{hungry}) - 5(1 - P(\text{hungry})) \\ = -10P(\text{hungry}) - 5 \end{aligned}$$



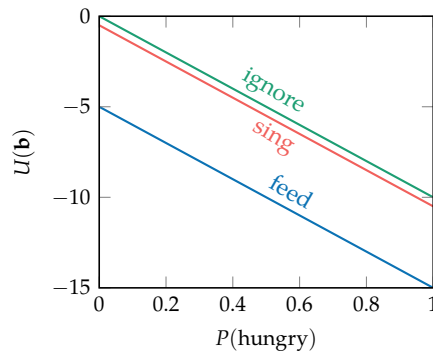
Singing to the baby yields an expected reward:

$$\begin{aligned} R(\text{hungry}, \text{sing})P(\text{hungry}) + R(\text{sated}, \text{sing})P(\text{sated}) \\ = -10.5P(\text{hungry}) - 0.5(1 - P(\text{hungry})) \\ = -10P(\text{hungry}) - 0.5 \end{aligned}$$

Ignoring the baby yields an expected reward:

$$\begin{aligned} R(\text{hungry}, \text{ignore})P(\text{hungry}) + R(\text{sated}, \text{ignore})P(\text{sated}) \\ = -10P(\text{hungry}) \end{aligned}$$

The expected reward for each action is plotted as follows over the belief space:



We find that under a one-step horizon, it is never optimal to feed or sing to the baby. The ignore action is dominant.

**Exercise 20.3.** Why does the implementation of value iteration in algorithm 20.8 call `expand` in algorithm 20.9 rather than evaluating the plan in algorithm 20.2 to obtain alpha vectors for each new conditional plan?

*Solution:* The plan evaluation method applies equation (20.8) recursively to evaluate the expected utility for a conditional plan. Conditional plans grow very large as the horizon increases. POMDP value iteration can save computation by using the alpha vectors for the subplans from the previous iteration:

$$U^\pi(s) = R(s, \pi()) + \gamma \left[ \sum_{s'} T(s' | s, \pi()) \sum_o O(o | \pi(), s') \alpha_{s'}^{\pi(o)} \right]$$

**Exercise 20.4.** Does the number of conditional plans increase faster with the number of actions or with the number of observations?

*Solution:* Recall that there are  $|\mathcal{A}|^{(|\mathcal{O}|^h - 1)/(|\mathcal{O}| - 1)}$  possible  $h$ -step plans. Exponential growth ( $n^x$ ) is faster than polynomial growth ( $x^n$ ), and we have better-than exponential growth in  $|\mathcal{O}|$  and polynomial growth in  $|\mathcal{A}|$ . The number of plans thus increases faster with respect to the number of observations. To demonstrate, let us use  $|\mathcal{A}| = 3$ ,  $|\mathcal{O}| = 3$ , and  $h = 3$  as a baseline. The baseline has 1,594,323 plans. Incrementing the number of actions results in 67,108,864 plans, whereas incrementing the number of observations results in 10,460,353,203 plans.

**Exercise 20.5.** Suppose that we have a patient and we are unsure whether they have a particular disease. We do have three diagnostic tests, each with different probabilities that they will correctly indicate whether the disease is present. While the patient is in our office, we have the option to administer multiple diagnostic tests in sequence. We observe the outcome of each diagnostic test immediately. In addition, we can repeat any diagnostic test multiple times, with the outcomes of all tests being conditionally independent of each other, given the presence or absence of the disease. When we are done with the tests, we decide whether to treat the disease or send the patient home without treatment. Explain how you would define the various components of a POMDP formulation.

*Solution:* We have three states:

1.  $s_{\text{no-disease}}$ : the patient does not have the disease
2.  $s_{\text{disease}}$ : the patient has the disease
3.  $s_{\text{terminal}}$ : the interaction is over (terminal state)

We have five actions:

1.  $a_1$ : administer test 1
2.  $a_2$ : administer test 2
3.  $a_3$ : administer test 3
4.  $a_{\text{treat}}$ : administer treatment and send patient home
5.  $a_{\text{stop}}$ : send patient home without treatment

We have three observations:

1.  $o_{\text{no-disease}}$ : the outcome of the test (if administered) indicates the patient does not have the disease
2.  $o_{\text{disease}}$ : the outcome of the test (if administered) indicates the patient has the disease
3.  $o_{\text{terminal}}$ : a test was not administered

The transition model would be deterministic, with

$$T(s' | s, a) = \begin{cases} 1 & \text{if } a \in \{a_{\text{treat}}, a_{\text{stop}}\} \wedge s' = s_{\text{terminal}} \\ 1 & \text{if } s = s' \\ 0 & \text{otherwise} \end{cases}$$

The reward function would be a function of the cost of administering treatment and each test, as well as the cost of not treating the disease if it is indeed present. The reward available from  $s_{\text{terminal}}$  is 0. The observation model assigns probabilities to correct and incorrect observations of the disease state as a result of a diagnostic test from one of the nonterminal states. The initial belief would assign our prior probability to whether the patient has the disease, with zero probability assigned to the terminal state.

**Exercise 20.6.** Why might we want to perform the same test multiple times in the previous exercise?

*Solution:* Depending on the probability of incorrect results, we may want to perform the same test multiple times to improve our confidence in whether the patient has the disease. The results of the tests are independent given the disease state.

**Exercise 20.7.** Suppose we have three alpha vectors,  $[1, 0]$ ,  $[0, 1]$ , and  $[\theta, \theta]$ , for a constant  $\theta$ . Under what conditions on  $\theta$  can we prune alpha vectors?

*Solution:* We can prune alpha vectors if  $\theta < 0.5$  or  $\theta > 1$ . If  $\theta < 0.5$ , then  $[\theta, \theta]$  is dominated by the other two alpha vectors. If  $\theta > 1$ , then  $[\theta, \theta]$  dominates the other two alpha vectors.

**Exercise 20.8.** We have  $\Gamma = \{[1, 0], [0, 1]\}$  and  $\alpha = [0.7, 0.7]$ . What belief  $\mathbf{b}$  maximizes the utility gap  $\delta$ , as defined by the linear program in equation (20.16)?

*Solution:* The alpha vectors in  $\Gamma$  are shown in blue and the alpha vector  $\alpha$  is shown in red. We care only about the region where  $0.3 \leq b_2 \leq 0.7$ , where  $\alpha$  dominates the alpha vectors in  $\Gamma$ ; in other words, where the red line is above the blue lines. The point where the gap between the red line and the maximum of the blue lines occurs at  $b_2 = 0.5$ , with a gap of  $\delta = 0.2$ . Hence, the belief that maximizes this gap is  $\mathbf{b} = [0.5, 0.5]$ .

