

27 Collaborative Agents

Many multiagent domains are collaborative, where all agents act independently in an environment while working toward a common shared objective. Applications range from robotic search and rescue to interplanetary exploration rovers. The *decentralized partially observable Markov decision process (Dec-POMDP)* captures the generality of POMGs while focusing on such collaborative agent settings.¹ The model is more amenable to scalable approximate algorithms because of its single shared objective, as opposed to finding an equilibrium among multiple individual agent objectives. This chapter presents the Dec-POMDP model, highlights its subclasses, and describes algorithms that solve them optimally and approximately.

¹D.S. Bernstein, R. Givan, N. Immerman, and S. Zilberstein, "The Complexity of Decentralized Control of Markov Decision Processes," *Mathematics of Operation Research*, vol. 27, no. 4, pp. 819–840, 2002. A more comprehensive introduction is provided by F.A. Oliehoek and C. Amato, *A Concise Introduction to Decentralized POMDPs*. Springer, 2016.

27.1 Decentralized Partially Observable Markov Decision Processes

A Dec-POMDP (algorithm 27.1) is a POMG where all agents share the same objective. Each agent $i \in \mathcal{I}$ selects a local action $a^i \in \mathcal{A}^i$ based on a history of local observations $o^i \in \mathcal{O}^i$. The true state of the system $s \in \mathcal{S}$ is shared by all agents. A single reward is generated by $R(s, \mathbf{a})$ based on state s and joint action \mathbf{a} . The goal of all agents is to maximize the shared expected reward over time under local partial observability. Example 27.1 describes a Dec-POMDP version of the predator-prey problem.

Consider a predator-prey hex world problem in which a team of predators \mathcal{I} strives to capture a single fleeing prey. The predators move independently. The prey moves randomly to a neighboring cell not occupied by a predator. The predators must work together to capture the prey.

Example 27.1. The collaborative predator-prey problem as a Dec-POMDP. Additional detail is provided in appendix F.15.

Many of the same challenges of POMGs persist in Dec-POMDPs, such as the general inability of agents to maintain a belief state. We focus on policies represented as conditional plans or controllers. The same algorithms introduced in the previous chapter can be used to evaluate policies. All that is required is to create a POMG with $R^i(s, \mathbf{a})$ for each agent i equal to the $R(s, \mathbf{a})$ from the Dec-POMDP.

```

struct DecPOMDP
   $\gamma$  # discount factor
   $\mathcal{I}$  # agents
   $\mathcal{S}$  # state space
   $\mathcal{A}$  # joint action space
   $\mathcal{O}$  # joint observation space
   $T$  # transition function
   $O$  # joint observation function
   $R$  # reward function
end

```

Algorithm 27.1. Data structure for a Dec-POMDP. The `joint` function from algorithm 24.2 allows the creation of all combinations of a set provided, such as \mathcal{A} or \mathcal{O} . The `tensorform` function converts the Dec-POMDP \mathcal{P} to a tensor representation.

27.2 Subclasses

There are many notable subclasses of Dec-POMDPs. Categorizing these subclasses is useful when designing algorithms that take advantage of their specific structure.

One attribute of interest is *joint full observability*, which is when each agent observes an aspect of the state, such that if they were to combine their observations, it would uniquely reveal the true state. The agents, however, do not share their observations. This property ensures that if $O(\mathbf{o} \mid \mathbf{a}, s') > 0$ then $P(s' \mid \mathbf{o}) = 1$. A Dec-POMDP with joint full observability is called a *decentralized Markov decision process (Dec-MDP)*. Both Dec-POMDP and Dec-MDP problems are *NEXP-complete* when the number of steps in the horizon is fewer than the number of states.²

In many settings, the state space of a Dec-POMDP is factored, one for each agent and one for the environment. This is called a *factored Dec-POMDP*. We have $\mathcal{S} = \mathcal{S}^0 \times \dots \times \mathcal{S}^k$, where \mathcal{S}^i is the factored state component associated with agent i and \mathcal{S}^0 is the factored state component associated with the general environment. For example, in the collaborative predator-prey problem, each agent has its own state factor for their location, and the position of the prey is associated with the environment component of the state space.

²In contrast with the complexity classes NP and PSPACE, it is known that NEXP is not in P. Hence, we can prove that Dec-MDPs and Dec-POMDPs do not allow for polynomial time algorithms. D. S. Bernstein, R. Givan, N. Immerman, and S. Zilberstein, "The Complexity of Decentralized Control of Markov Decision Processes," *Mathematics of Operation Research*, vol. 27, no. 4, pp. 819–840, 2002.

In some problems, a factored Dec-POMDP may have one or more of the following properties:

- *Transition independence*, where agents may not affect each other's state:

$$T(\mathbf{s}' | \mathbf{s}, \mathbf{a}) = T^0(s^{0'} | s^0) \prod_i T^i(s^{i'} | s^i, a^i) \quad (27.1)$$

- *Observation independence*, where the observations of agents depend only on their local state and actions:

$$O(\mathbf{o} | \mathbf{a}, \mathbf{s}') = \prod_i O^i(o^i | a^i, s^{i'}) \quad (27.2)$$

- *Reward independence*, where the reward can be decomposed into multiple independent pieces:³

$$R(\mathbf{s}, \mathbf{a}) = R^0(s^0) + \sum_i R^i(s^i, a^i) \quad (27.3)$$

³ Here, we show the combination of the reward components as a summation, but any monotonically nondecreasing function can be used instead and preserve reward independence.

The computational complexity can vary significantly depending on which of these independence properties are satisfied, as summarized in table 27.1. It is important to take these independences into account when modeling a problem to improve scalability.

Independence	Complexity
Transitions, observations, and rewards	P-complete
Transitions and observations	NP-complete
Any other subset	NEXP-complete

Table 27.1. The complexity of factored Dec-POMDPs with different independence assumptions.

A *network distributed partially observable Markov decision process (ND-POMDP)* is a Dec-POMDP with transition and observation independence and a special reward structure. The reward structure is represented by a coordination graph. In contrast with the graphs used earlier in this book, a *coordination graph* is a type of hypergraph, which allows edges to connect any number of nodes. The nodes in the ND-POMDP hypergraph correspond to the various agents. The edges relate to interactions between the agents in the reward function. An ND-POMDP associates with each edge j in the hypergraph a reward component R_j that depends on the state and action components to which the edge connects. The reward function in

an ND-POMDP is simply the sum of the reward components associated with the edges. Figure 27.1 shows a coordination graph resulting in a reward function that can be decomposed as follows:

$$R_{123}(s_1, s_2, s_3, a_1, a_2, a_3) + R_{34}(s_3, s_4, a_3, a_4) + R_5(s_5, a_5) \quad (27.4)$$

Sensor network and target tracking problems are often framed as ND-POMDPs.

The ND-POMDP model is similar to the transition and observation independent Dec-MDP model, but it does not make the joint full observability assumption. Even if all observations are shared, the true state of the world may not be known. Furthermore, even with factored transitions and observations, a policy in an ND-POMDP is a mapping from observation histories to actions, unlike the transition and observation Dec-MDP case, in which policies are mappings from local states to actions. The worst-case complexity remains the same as for a Dec-POMDP, but algorithms for ND-POMDPs are typically much more scalable in the number of agents. Scalability can increase as the coordination graph becomes less connected.

If the agents are able to communicate their actions and observations perfectly without penalty, then they are able to maintain a collective belief state. This model is called a *multiagent MDP (MMDP)* or a *multiagent POMDP (MPOMDP)*. MMDPs and MPOMDPs can also result when there is transition, observation, and reward independence. Any MDP or POMDP algorithm discussed in earlier chapters can be applied to solve these problems.

Table 27.2 summarizes some of these subclasses. Figure 27.2 illustrates the relationships among the models discussed in this book.

Agents	Observability	Communication	Model
Single	Full	—	MDP
Single	Partial	—	POMDP
Multiple	Full	Free	MMDP
Multiple	Full	General	MMDP
Multiple	Joint full	Free	MMDP
Multiple	Joint full	General	Dec-MDP
Multiple	Partial	Free	MPOMDP
Multiple	Partial	General	Dec-POMDP

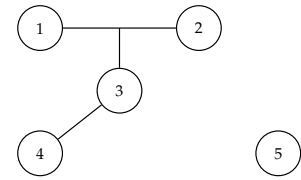


Figure 27.1. An ND-POMDP structure with five agents. There are three hyperedges: one involving agents 1, 2, and 3; another involving agents 3 and 4; and another involving agent 5 on its own.

Table 27.2. Dec-POMDP subclasses categorized by type and computational complexity. “Observability” refers to the degree to which the shared state is observable. “Communication” refers to whether the cooperative agents can freely share all observations with each other. Free communication happens outside the model (e.g., a high-speed wireless connection in robots). General communication is when agents do not have this available and must communicate (typically imperfectly) via their actions.

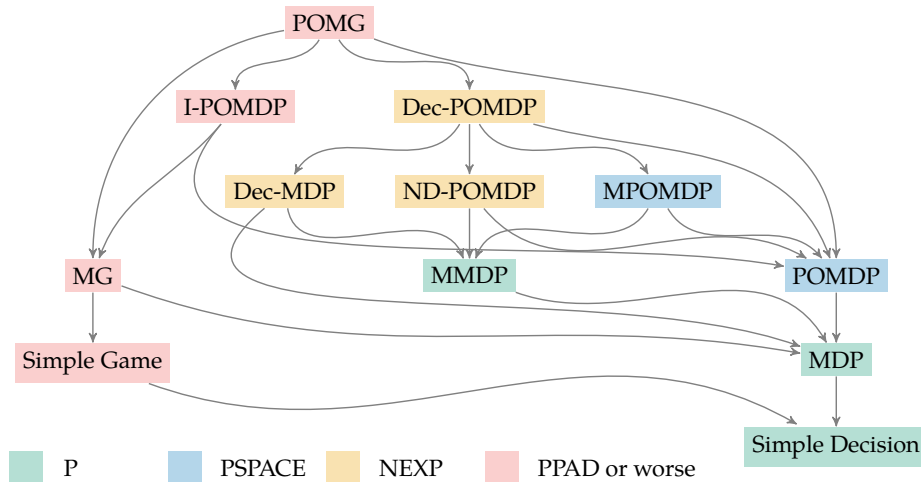


Figure 27.2. A taxonomy for the models discussed in this book. Parents generalize their children in this diagram. For example, Dec-POMDPs generalize POMDPs by supporting multiple agents. The color of the nodes indicate computational complexity, as indicated in the key at the bottom left of the figure. The complexities listed here are for the common model, policy, and objective formulations presented in this book. For a more detailed treatment, see C. Papadimitriou and J. Tsitsiklis, “The Complexity of Markov Decision Processes,” *Mathematics of Operation Research*, vol. 12, no. 3, pp. 441–450, 1987. Also, see S. Seuken and S. Zilberstein, “Formal Models and Algorithms for Decentralized Decision Making Under Uncertainty,” *Autonomous Agents and Multi-Agent Systems*, vol. 17, no. 2, pp. 190–250, 2008.

27.3 Dynamic Programming

The *dynamic programming* algorithm for Dec-POMDPs applies the Bellman backup at each step and prunes dominated policies. This process is identical to dynamic programming for POMGs except that each agent shares the same reward. Algorithm 27.2 implements this procedure.

```

struct DecPOMDPDynamicProgramming
  b # initial belief
  d # depth of conditional plans
end

function solve(M::DecPOMDPDynamicProgramming, P::DecPOMDP)
  T, S, A, O, T, O, R, γ = P.T, P.S, P.A, P.O, P.T, P.O, P.R, P.γ
  R'(s, a) = [R(s, a) for i in I]
  P' = POMG(γ, T, S, A, O, T, O, R')
  M' = POMGDynamicProgramming(M.b, M.d)
  return solve(M', P')
end

```

Algorithm 27.2. Dynamic programming computes the optimal joint policy π for a Dec-POMDP \mathcal{P} , given an initial belief \mathbf{b} and horizon depth \mathbf{d} . We can directly use the POMG algorithm, as Dec-POMDPs are a special collaborative class of POMGs.

27.4 Iterated Best Response

Instead of exploring joint policies directly, we can perform a form of *iterated best response* (algorithm 27.3). In this approach, we iteratively select an agent and compute a best response policy, assuming that the other agents are following a fixed policy.⁴ This approximate algorithm is typically fast because it is choosing the best policy for only one agent at a time. Moreover, since all agents share the same reward, it tends to terminate after relatively few iterations.

Iterated best response begins with a random initial joint policy π_1 . The process randomly iterates over the agents. If agent i is selected, its policy π^i is updated with a best response to the other agents' fixed policies π^{-i} with initial belief distribution b :

$$\pi^i \leftarrow \arg \max_{\pi^{i'}} U^{\pi^{i'}, \pi^{-i}}(b) \quad (27.5)$$

with ties favoring the current policy. This process can terminate when agents stop changing their policies.

While this algorithm is fast and guaranteed to converge, it does not always find the best joint policy. It relies on iterated best response to find a Nash equilibrium, but there may be many Nash equilibria, with different utilities associated with them. This approach will find only one of them.

27.5 Heuristic Search

Instead of expanding all joint policies, *heuristic search* (algorithm 27.4) explores a fixed number of policies,⁵ which, stored over iterations, prevents exponential growth. The heuristic exploration guides the search by attempting to expand the best joint policies only until depth d is reached.

Each iteration k of the algorithm keeps a set of joint policies Π_k . This set initially consists of all one-step conditional plans. Subsequent iterations begin by fully expanding the conditional plans. The goal is to add a fixed number of these for the next iteration.

We prioritize the policies that are more likely to maximize the utility when deciding among the conditional plans to add to the set. However, since we expand the conditional plans from the bottom up, we cannot simply evaluate the policies from the initial belief state b . Instead, we need an estimate of the belief $d - k$ steps into the future, which we compute by taking random actions and simulating state

⁴This type of algorithm is also called *joint equilibrium-based search for policies (JESP)*. R. Nair, M. Tambe, M. Yokoo, D. Pynadath, and S. Marsella, "Taming Decentralized POMDPs: Towards Efficient Policy Computation for Multiagent Settings," in *International Joint Conference on Artificial Intelligence (IJCAI)*, 2003. It can be improved further by performing dynamic programming.

⁵This approach is also known as *memory-bounded dynamic programming (MBDP)*. S. Seuken and S. Zilberstein, "Memory-Bounded Dynamic Programming for Dec-POMDPs," in *International Joint Conference on Artificial Intelligence (IJCAI)*, 2007. There are other heuristic search algorithms as well, such as *multiagent A* (MMA*)*. D. Szer, F. Charpillet, and S. Zilberstein, "MAA*: A Heuristic Search Algorithm for Solving Decentralized POMDPs," in *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2005.

```

struct DecPOMDPIteratedBestResponse
    b      # initial belief
    d      # depth of conditional plans
    k_max  # number of iterations
end

function solve(M::DecPOMDPIteratedBestResponse, P::DecPOMDP)
    I, S, A, O, T, O, R, γ = P.I, P.S, P.A, P.O, P.T, P.O, P.R, P.γ
    b, d, k_max = M.b, M.d, M.k_max
    R'(s, a) = [R(s, a) for i in I]
    P' = POMG(γ, I, S, A, O, T, O, R')
    Π = create_conditional_plans(P, d)
    π = [rand(Π[i]) for i in I]
    for k in 1:k_max
        for i in shuffle(I)
            π'(π) = Tuple(j == i ? π : π[j] for j in I)
            Ui(π) = utility(P', b, π'(π))[i]
            π[i] = argmax(Ui, Π[i])
        end
    end
    return Tuple(π)
end

```

Algorithm 27.3. Iterated best response for a collaborative Dec-POMDP \mathcal{P} performs a deterministic best response for each agent to rapidly search the space of conditional plan policies. The `solve` function executes this procedure for up to `k_max` steps, maximizing the value at an initial belief `b` for conditional plans of depth `d`.

transitions and observations, updating the belief along the way. This belief at iteration k is denoted as b_k . For each available joint policy $\pi \in \Pi_k$, the utility $U^\pi(b_k)$ is examined to find a utility-maximizing joint policy to add. Example 27.2 demonstrates the process.

27.6 Nonlinear Programming

We can use *nonlinear programming* (NLP) (algorithm 27.5) to find an optimal joint controller policy representation of a fixed size.⁶ This method generalizes the NLP approach for POMDPs from section 23.3.

⁶C. Amato, D.S. Bernstein, and S. Zilberstein, “Optimizing Fixed-Size Stochastic Controllers for POMDPs and Decentralized POMDPs,” *Autonomous Agents and Multi-Agent Systems*, vol. 21, no. 3, pp. 293–320, 2010.

```

struct DecPOMDPHeuristicSearch
    b # initial belief
    d # depth of conditional plans
    π_max # number of policies
end

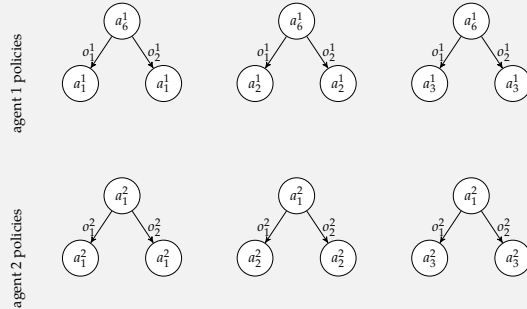
function solve(M::DecPOMDPHeuristicSearch, P::DecPOMDP)
    T, S, A, O, T, O, R, γ = P.T, P.S, P.A, P.O, P.T, P.O, P.R, P.γ
    b, d, π_max = M.b, M.d, M.π_max
    R'(s, a) = [R(s, a) for i in T]
    P' = POMG(γ, T, S, A, O, T, O, R')
    Π = [[ConditionalPlan(ai) for ai in A[i]] for i in T]
    for t in 1:d
        allΠ = expand_conditional_plans(P, Π)
        Π = [[] for i in T]
        for z in 1:π_max
            b' = explore(M, P, t)
            π = argmax(π → first(utility(P', b', π)), joint(allΠ))
            for i in T
                push!(Π[i], π[i])
                filter!(πi → πi != π[i], allΠ[i])
            end
        end
    end
    return argmax(π → first(utility(P', b, π)), joint(Π))
end

function explore(M::DecPOMDPHeuristicSearch, P::DecPOMDP, t)
    T, S, A, O, T, O, R, γ = P.T, P.S, P.A, P.O, P.T, P.O, P.R, P.γ
    b = copy(M.b)
    b' = similar(b)
    s = rand(SetCategorical(S, b))
    for τ in 1:t
        a = Tuple(rand(Ai) for Ai in A)
        s' = rand(SetCategorical(S, [T(s,a,s') for s' in S]))
        o = rand(SetCategorical(joint(O), [O(a,s',o) for o in joint(O)]))
        for (i', s') in enumerate(S)
            po = O(a, s', o)
            b'[i'] = po*sum(T(s,a,s')*b[i] for (i,s) in enumerate(S))
        end
        normalize!(b', 1)
        b, s = b', s'
    end
    return b'
end

```

Algorithm 27.4. Memory-bounded heuristic search uses a heuristic function to search the space of conditional plans for a Dec-POMDP \mathcal{P} . The `solve` function tries to maximize the value at an initial belief \mathbf{b} for joint conditional plans of depth d . The `explore` function generates a belief \mathbf{t} steps into the future by taking random actions and simulating actions and observations. The algorithm is memory-bounded, keeping only π_{\max} conditional plans per agent.

Consider the collaborative predator-prey problem shown at right. We apply heuristic search to a depth of $d = 3$, with three policies retained at each iteration. After iteration $k = 1$, the policies are



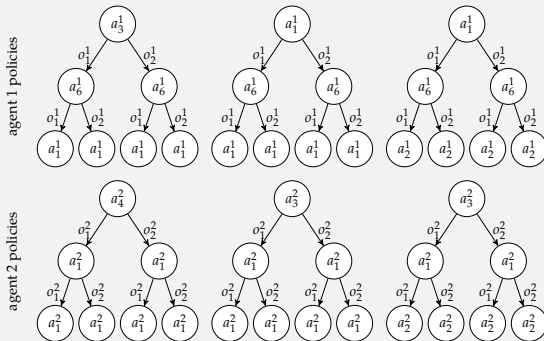
Example 27.2. Heuristic search exploration and conditional plan expansion for the collaborative predator-prey hex world problem shown here. The predators are red and green. The prey is blue.



At the next iteration $k = 2$, heuristic search again starts at the initial belief and takes $d - k = 3 - 2 = 1$ steps following the heuristic exploration. The explored beliefs used to select the next three conditional plans are

$$\begin{array}{lll}
 b_1 = [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.17 & b_2 = [0.0, 0.21, 0.03, 0.0, 0.04, 0.01, 0.0 & b_3 = [0.0, 0.03, 0.01, 0.0, 0.03, 0.01, 0.0 \\
 0.0, 0.03, 0.01, 0.0, 0.0, 0.05, 0.0 & 0.05, 0.01, 0.0, 0.08, 0.03, 0.0, 0.0 & 0.15, 0.05, 0.0, 0.01, 0.0, 0.0, 0.0 \\
 0.01, 0.23, 0.0, 0.08, 0.01, 0.0, 0.0 & 0.01, 0.0, 0.0, 0.01, 0.08, 0.34, 0.03 & 0.0, 0.0, 0.0, 0.03, 0.06, 0.11, 0.32 \\
 0.14, 0.0, 0.03, 0.22, 0.0, 0.01] & 0.02, 0.05, 0.01, 0.0, 0.01, 0.0] & 0.06, 0.03, 0.01, 0.01, 0.04, 0.06]
 \end{array}$$

The policies after iteration $k = 2$ are



The beliefs were used to determine the root node's action and the two subtrees below it. These subtrees are built from the prior iteration's trees.

Given a fixed set of nodes X^i for each agent i , initial belief b , and initial joint nodes \mathbf{x}_1 , the optimization problem is

$$\begin{aligned}
 & \underset{U, \psi, \eta}{\text{maximize}} && \sum_s b(s) U(\mathbf{x}_1, s) \\
 & \text{subject to} && U(\mathbf{x}, s) = \sum_{\mathbf{a}} \prod_i \psi^i(a^i | x^i) \left(R(s, \mathbf{a}) + \gamma \sum_{s'} T(s' | s, \mathbf{a}) \sum_{\mathbf{o}} O(\mathbf{o} | \mathbf{a}, s') \sum_{\mathbf{x}'} \prod_i \eta^i(x^{i'} | x^i, a^i, o^i) U(\mathbf{x}', s') \right) \\
 & && \text{for all } \mathbf{x}, s \\
 & && \psi^i(a^i | x^i) \geq 0 \quad \text{for all } i, x^i, a^i \\
 & && \sum_a \psi^i(a^i | x^i) = 1 \quad \text{for all } i, x^i \\
 & && \eta^i(x^{i'} | x^i, a^i, o^i) \geq 0 \quad \text{for all } i, x^i, a^i, o^i, x^{i'} \\
 & && \sum_{x^{i'}} \eta^i(x^{i'} | x^i, a^i, o^i) = 1 \quad \text{for all } i, x^i, a^i, o^i
 \end{aligned} \tag{27.6}$$

27.7 Summary

- Dec-POMDPs are fully cooperative POMGs that model a team of agents working together toward a shared goal, each acting individually using only local information.
- Because determining a belief state is infeasible, as in POMGs, policies are generally represented as conditional plans or controllers, allowing each agent to map individual sequences of observations to individual actions.
- Many subclasses of Dec-POMDPs exist, with different degrees of computational complexity.
- Dynamic programming computes the value function iteratively, pruning dominated policies as it iterates using a linear program.
- Iterated best response computes a best utility-maximizing response policy for a single agent at a time, iteratively converging to a joint equilibrium.
- Heuristic search searches a fixed subset of policies at each iteration, guided by a heuristic.
- Nonlinear programming can be used to generate controllers of a fixed size.

```

struct DecPOMDPNonlinearProgramming
    b # initial belief
    ℓ # number of nodes for each agent
end

function tensorform( $\mathcal{P}$ ::DecPOMDP)
     $\mathcal{T}, S, \mathcal{A}, \mathcal{O}, R, T, O = \mathcal{P}.\mathcal{T}, \mathcal{P}.S, \mathcal{P}.\mathcal{A}, \mathcal{P}.\mathcal{O}, \mathcal{P}.R, \mathcal{P}.T, \mathcal{P}.\mathcal{O}$ 
     $\mathcal{T}' = \text{eachindex}(\mathcal{T})$ 
     $S' = \text{eachindex}(S)$ 
     $\mathcal{A}' = [\text{eachindex}(\mathcal{A}_i) \text{ for } \mathcal{A}_i \text{ in } \mathcal{A}]$ 
     $\mathcal{O}' = [\text{eachindex}(\mathcal{O}_i) \text{ for } \mathcal{O}_i \text{ in } \mathcal{O}]$ 
     $R' = [R(s,a) \text{ for } s \text{ in } S, a \text{ in } \text{joint}(\mathcal{A})]$ 
     $T' = [T(s,a,s') \text{ for } s \text{ in } S, a \text{ in } \text{joint}(\mathcal{A}), s' \text{ in } S]$ 
     $O' = [O(a,s',o) \text{ for } a \text{ in } \text{joint}(\mathcal{A}), s' \text{ in } S, o \text{ in } \text{joint}(\mathcal{O})]$ 
    return  $\mathcal{T}', S', \mathcal{A}', \mathcal{O}', R', T', O'$ 
end

function solve( $M$ ::DecPOMDPNonlinearProgramming,  $\mathcal{P}$ ::DecPOMDP)
     $\mathcal{P}, \gamma, b = \mathcal{P}, \mathcal{P}.\gamma, M.b$ 
     $\mathcal{T}, S, \mathcal{A}, \mathcal{O}, R, T, O = \text{tensorform}(\mathcal{P})$ 
     $X = [\text{collect}(1:M.\ell) \text{ for } i \text{ in } \mathcal{T}]$ 
     $\text{jointX}, \text{jointA}, \text{jointO} = \text{joint}(X), \text{joint}(\mathcal{A}), \text{joint}(\mathcal{O})$ 
     $x1 = \text{jointX}[1]$ 
     $\text{model} = \text{Model}(\text{Ipopt}.\text{Optimizer})$ 
    @variable( $\text{model}$ ,  $U[\text{jointX}, S]$ )
    @variable( $\text{model}$ ,  $\psi[i=\mathcal{T}, X[i], \mathcal{A}[i]] \geq 0$ )
    @variable( $\text{model}$ ,  $\eta[i=\mathcal{T}, X[i], \mathcal{A}[i], \mathcal{O}[i], X[i]] \geq 0$ )
    @objective( $\text{model}$ , Max,  $b \cdot U[x1, :]$ )
    @NLconstraint( $\text{model}$ , [ $x=\text{jointX}, s=S$ ],
         $U[x, s] == (\text{sum}(\text{prod}(\psi[i, x[i], a[i]] \text{ for } i \text{ in } \mathcal{T})$ 
             $\ast (R[s, y] + \gamma \ast \text{sum}(T[s, y, s'] \ast \text{sum}(O[y, s', z]$ 
                 $\ast \text{sum}(\text{prod}(\eta[i, x[i], a[i], o[i], x'[i]] \text{ for } i \text{ in } \mathcal{T})$ 
                     $\ast U[x', s'] \text{ for } x' \text{ in } \text{jointX})$ 
                for  $(z, o) \text{ in } \text{enumerate}(\text{jointO})$  for  $s' \text{ in } S)$ 
                for  $(y, a) \text{ in } \text{enumerate}(\text{jointA}))$ 
             $\ast U[x, s]$ 
         $\text{sum}(\psi[i, xi, ai] \text{ for } ai \text{ in } \mathcal{A}[i]) == 1$ 
         $\text{sum}(\eta[i, xi, ai, oi, xi'] \text{ for } xi' \text{ in } X[i]) == 1$ 
     $\text{optimize!}(\text{model})$ 
     $\psi', \eta' = \text{value}(\psi), \text{value}(\eta)$ 
    return  $[\text{ControllerPolicy}(\mathcal{P}, X[i],$ 
         $\text{Dict}((xi, \mathcal{P}.\mathcal{A}[i][ai]) \Rightarrow \psi'[i, xi, ai]$ 
            for  $xi \text{ in } X[i], ai \text{ in } \mathcal{A}[i]),$ 
         $\text{Dict}((xi, \mathcal{P}.\mathcal{A}[i][ai], \mathcal{P}.\mathcal{O}[i][oi], xi') \Rightarrow \eta'[i, xi, ai, oi, xi']$ 
            for  $xi \text{ in } X[i], ai \text{ in } \mathcal{A}[i], oi \text{ in } \mathcal{O}[i], xi' \text{ in } X[i]))$ 
    for  $i \text{ in } \mathcal{T}]$ 
end

```

Algorithm 27.5. NLP computes the optimal joint controller policy π for a Dec-POMDP \mathcal{P} , given an initial belief b and number of controller nodes ℓ for each agent. This generalizes the NLP solution in algorithm 23.5.

27.8 Exercises

Exercise 27.1. Why is a Dec-MDP with joint full observability different from agents knowing the state?

Solution: Full joint observability means if agents were to share their individual observations, then the team would know the true state. This can be done offline during planning. Thus in Dec-MDPs, the true state is essentially known during planning. The issue is that it requires agents to share their individual observations, which cannot be done online during execution. Therefore, planning still needs to reason about the uncertain observations made by the other agents.

Exercise 27.2. Propose a fast algorithm for a Dec-MDP with transition, observation, and reward independence. Prove that it is correct.

Solution: If a factored Dec-MDP satisfies all three independence assumptions, then we can solve it as $|\mathcal{I}|$ separate MDPs. The resulting policy π^i for each agent i 's MDP can then be combined to form the optimal joint policy. To prove this fact, consider the utility of each agent's individual MDP:

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

As in equation (26.1), $\pi^i()$ refers to the root action of i 's conditional plan, and $\pi^i(o^i)$ refers to i 's subplans after making observation o^i . We sum up each of their individual contributions as follows:

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

We can combine T^i and O^i into a single probability distribution P , move the summation, and apply the definition of reward independence:

$$\begin{aligned} \sum_i U^{\pi^i}(s) &= \sum_i \left[R(s^i, \pi^i()) + \gamma \left[\sum_{s^{i'}} P(s^{i'} | s^i, \pi^i()) \sum_{o^i} P(o^i | \pi^i(), s^{i'}) U^{\pi^i(o^i)}(s^{i'}) \right] \right] \\ &= \sum_i R(s^i, \pi^i()) + \sum_i \left[\gamma \left[\sum_{s^{i'}} P(s^{i'} | s^i, \pi^i()) \sum_{o^i} P(o^i | \pi^i(), s^{i'}) U^{\pi^i(o^i)}(s^{i'}) \right] \right] \\ &= R(s, \pi()) + \sum_i \left[\gamma \left[\sum_{s^{i'}} P(s^{i'} | s^i, \pi^i()) \sum_{o^i} P(o^i | \pi^i(), s^{i'}) U^{\pi^i(o^i)}(s^{i'}) \right] \right] \end{aligned}$$

Now, we marginalize over all successors s and observations o . Because of the transition and observation independence, we can freely condition the distributions on these other non- i state and observation factors, which is the same as conditioning on s and o . We can then apply the definition of transition and observation independence. Finally, we can move the summation in and recognize $U^\pi(s)$ results:

$$\begin{aligned}
\sum_i U^{\pi^i}(s) &= R(s, \pi()) + \sum_i \left[\gamma \left[\sum_{s'} P(s' | s^i, \pi^i()) \sum_o P(o | \pi^i(), s^{i'}) U^{\pi^i(o^i)}(s^{i'}) \right] \right] \\
&= R(s, \pi()) + \sum_i \left[\gamma \left[\sum_{s'} P(s^{0'} | s^0) \prod_j P(s^{j'} | s^i, \pi^i()) \sum_o \prod_j P(o^j | \pi^i(), s^{i'}) U^{\pi^i(o^i)}(s^{i'}) \right] \right] \\
&= R(s, \pi()) + \sum_i \left[\gamma \left[\sum_{s'} P(s^{0'} | s^0) \prod_j P(s^{j'} | s, \pi()) \sum_o \prod_j P(o^j | \pi(), s') U^{\pi^i(o^i)}(s^{i'}) \right] \right] \\
&= R(s, \pi()) + \sum_i \left[\gamma \left[\sum_{s'} T(s' | s, \pi()) \sum_o O(o | \pi(), s') U^{\pi^i(o^i)}(s^{i'}) \right] \right] \\
&= R(s, \pi()) + \gamma \left[\sum_{s'} T(s' | s, \pi()) \sum_o O(o | \pi(), s') \left[\sum_i U^{\pi^i(o^i)}(s^{i'}) \right] \right] \\
&= R(s, \pi()) + \gamma \left[\sum_{s'} T(s' | s, \pi()) \sum_o O(o | \pi(), s') U^{\pi(o)}(s') \right] \\
&= U^\pi(s)
\end{aligned}$$

This is the Dec-MDP utility function derived from equation (26.1), completing the proof.

Exercise 27.3. How can we use an MMDP or MPOMDP as a heuristic in Dec-POMDP heuristic search?

Solution: We can assume free communication for planning. At each time step t , all agents know \mathbf{a}_t and \mathbf{o}_t , allowing us to maintain a multiagent belief b_t , resulting in an MPOMDP. This MPOMDP solution can be used as a heuristic to guide the search of policy trees. Alternatively, we create a heuristic where we assume that the true state and joint actions are known. This results in an MMDP, and it can also be used as a heuristic. These assumptions are used only for planning. Execution is still a Dec-POMDP wherein agents receive individual observations without free communication. Either heuristic results in a joint policy $\hat{\pi}$ for heuristic exploration.

Exercise 27.4. How can we compute a best response controller? Describe how this could be used in an iterated best response.

Solution: For an agent i , the best response controller X^i , ψ^i , and η^i can be computed by solving a nonlinear program. The program is similar to what is given in section 27.6, except that X^{-i} , ψ^{-i} , and η^{-i} are now given and are no longer variables:

$$\begin{aligned} & \text{maximize}_{U, \psi^i, \eta^i} \sum_s b(s) U(\mathbf{x}_1, s) \\ & \text{subject to } U(\mathbf{x}, s) = \sum_{\mathbf{a}} \prod_i \psi^i(a^i | x^i) \left(R(s, \mathbf{a}) + \gamma \sum_{s'} T(s' | s, \mathbf{a}) \sum_{\mathbf{o}} O(\mathbf{o} | \mathbf{a}, s') \sum_{\mathbf{x}'} \prod_i \eta^i(x^{i'} | x^i, a^i, o^i) U(\mathbf{x}', s') \right) \\ & \hspace{20em} \text{for all } \mathbf{x}, s \\ & \psi^i(a^i | x^i) \geq 0 \quad \text{for all } x^i, a^i \\ & \sum_a \psi^i(a^i | x^i) = 1 \quad \text{for all } x^i \\ & \eta^i(x^{i'} | x^i, a^i, o^i) \geq 0 \quad \text{for all } x^i, a^i, o^i, x^{i'} \\ & \sum_{x^{i'}} \eta^i(x^{i'} | x^i, a^i, o^i) = 1 \quad \text{for all } x^i, a^i, o^i \end{aligned}$$

Adapting algorithm 27.3 for controller policies, this program replaces the inner best response operation.