# Decision Making under Uncertainty
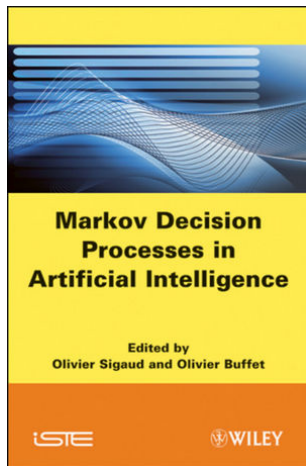## MDPs and POMDPs

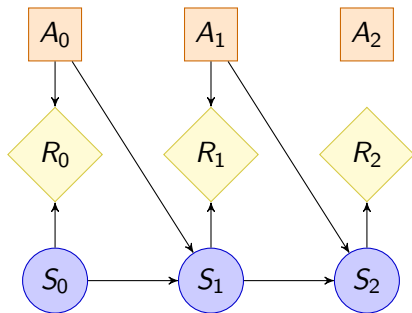Mykel J. Kochenderfer

27 January 2014

# Recommended reference

### Markov Decision Processes in Artificial Intelligence
edited by Sigaud and Buffet
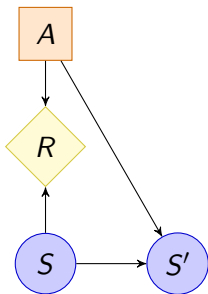


- Surveys a broad range of related topics
- Published in 2008, contains many recent references
- Contains psuedocode for algorithms
- More rigorous treatments can be found elsewhere

# Stochastic problems



- Agent chooses action $A_t$ at time $t$ based on observing state $S_t$
- State evolves probabilistically based on current state and action taken by agent (Markov assumption)
- Objective is to maximize sum of rewards $R_1, \ldots$

# Markov decision process



- If transition and utility models are stationary (does not vary with time), then a dynamic Bayesian network can be used
- Problem is known as a Markov decision process (MDPs)
- Defined by transition model $T(s' \mid s, a)$
  and reward model $R(s, a)$ (often assumed deterministic)

# Example Markov decision process

Aircraft collision avoidance



intruder aircraft

own aircraft

- ▶ Own aircraft must choose to stay level, climb, or descend
- ▶ At each step, $-1$ for collision, $-0.01$ for climb or descend
- ▶ State determined by altitude, closure rate, and vertical rates
- ▶ Intruder aircraft flies around randomly

Optimal behavior determined by reward and transition model

# Models of optimal behavior

- In the finite-horizon model agent should optimize expected reward for the next $H$ steps: $E(\sum_{t=0}^{H} r_t)$
- Continuously executing $H$-step optimal actions is known as receding horizon control
- In the infinite-horizon discounted model agent should optimize $E(\sum_{t=0}^{\infty} \gamma^t r_t)$
- Discount factor $0 \leq \gamma < 1$ can be thought of as an interest rate (reward now is worth more than reward in the future)
- Discounting keeps utility of an infinite sequence finite

# Outline

# Policies and utilities of states

- A policy $\pi$ specifies what action to execute from every possible state
- Action to execute from state $s$ according to $\pi$ is denoted $\pi(s)$
- Expected utility of executing $\pi$ when starting from $s$ is denoted $U^\pi(s)$
- Optimal policy $\pi^*$ is one that maximizes expected utility: $\pi^*(s) = \arg\max_\pi U^\pi(s)$

# Iterative policy evaluation

- Incrementally compute expected utility after $k$ steps of executing $\pi$
- $U_0^\pi(s) = 0$
- $U_1^\pi(s) = R(s, \pi(s))$
- ...
- $U_k^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s' \mid s, \pi(s)) U_{k-1}^\pi(s')$

This kind of iterative calculation is called dynamic programming

# Policy evaluation

- For an infinite horizon,
  $U^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s' \mid s, \pi(s)) U^\pi(s')$
- Can compute this arbitrarily well with enough iterations of iterative policy evaluation
- Alternative is to just solve system of $N$ linear equations, where $N$ is the number of states, requiring $O(n^3)$ time
- $U^\pi = R^\pi + \gamma T^\pi U^\pi$ (in matrix form)
- $U^\pi - \gamma T^\pi U^\pi = R^\pi$
- $(I - \gamma T^\pi) U^\pi = R^\pi$
- $U^\pi = (I - \gamma T^\pi)^{-1} R^\pi$
  (in Matlab, best to use matrix left division)

# Policy iteration

- Policy iteration is one way to compute an optimal policy $\pi^*$
- The algorithm starts with any policy $\pi_0$ and iterates the following steps
    1. Policy evaluation: given $\pi_i$ compute $U^{\pi_i}$
    2. Policy improvement: compute new policy from $U^{\pi_i}$
       $\pi_{i+1}(s) = \arg\max_a [R(s, a) + \gamma \sum_{s'} T(s' \mid s, a) U^{\pi_i}(s')]$
- Algorithm terminates when there is no more improvement
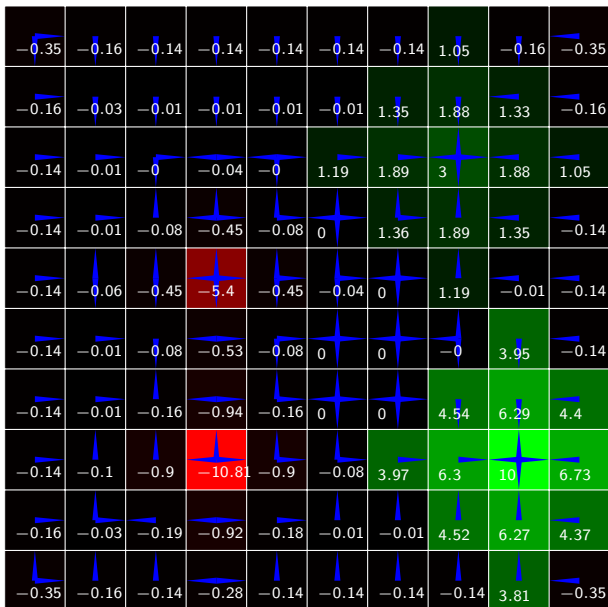- Since every step leads to improvement and there are finitely many policies, algorithm terminates at optimal solution

# Value iteration

- An alternative algorithm is value iteration
- Bellman equation says value of optimal policy is given by
  $U^*(s) = \max_a[R(s, a) + \gamma \sum_{s'} T(s' \mid s, a)U^*(s')]$
- $U_0^*(s) = 0$ (initialization)
- $U_1^*(s) = \max_a[R(s, a) + \gamma \sum_{s'} T(s' \mid s, a)U_0^*(s')]$
- ...
- $U_k^*(s) = \max_a[R(s, a) + \gamma \sum_{s'} T(s' \mid s, a)U_{k-1}^*(s')]$
- $U_k^* \to U^*$ as $k \to \infty$
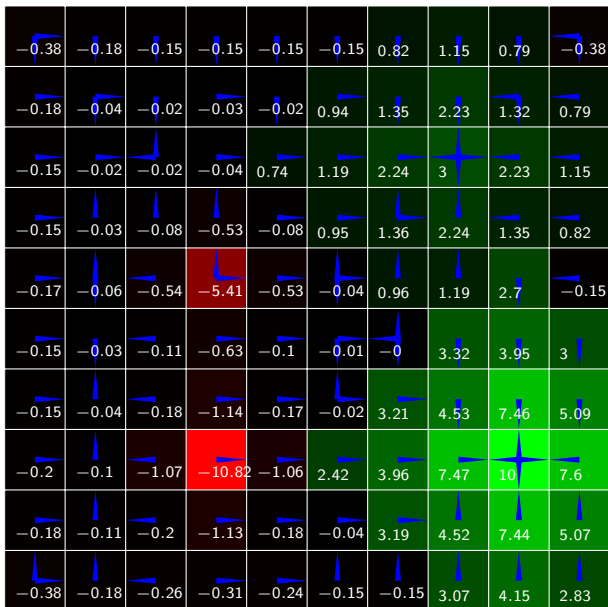- $\pi(s) = \arg\max_a[R(s, a) + \gamma \sum_{s'} T(s' \mid s, a)U(s')]$

# Example grid world

- States: cells in $10 \times 10$ grid
- Actions: up, down, left, and right
- Transition model: 0.7 chance of moving in intended direction, uniform in other directions
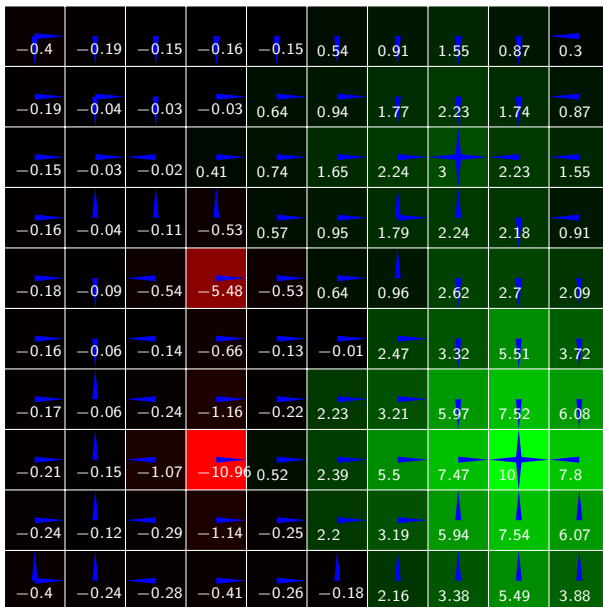- Reward: two states have cost (and are not terminal), two have reward (and are terminal), $-1$ for wall crash
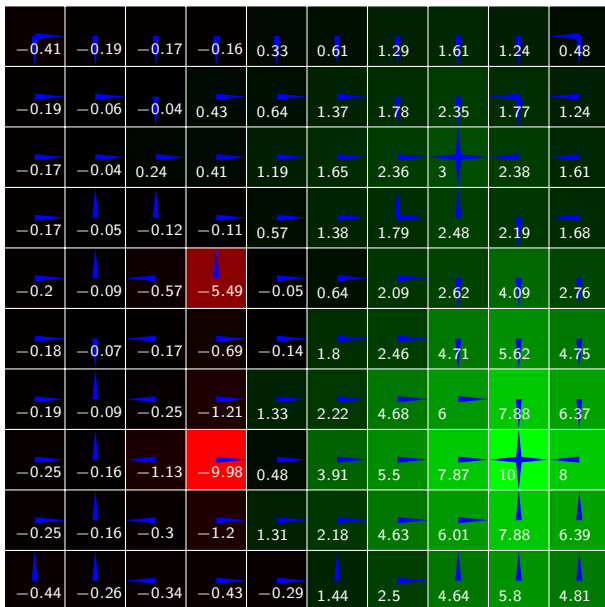
# Grid world: Iteration 1 ($\gamma = 0.9$)

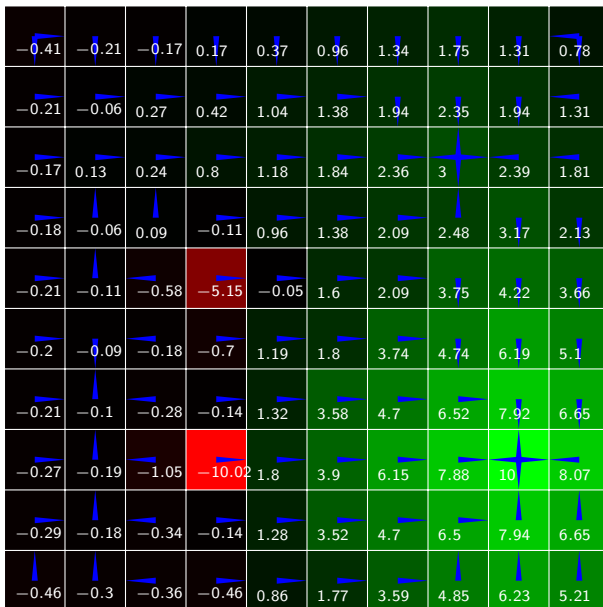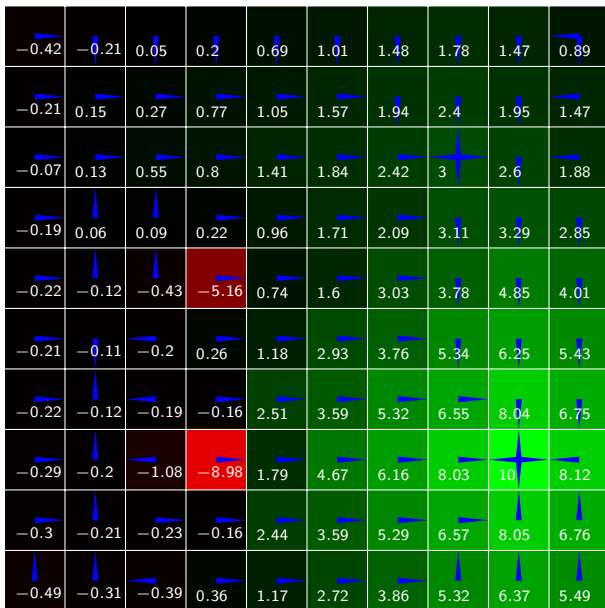# Grid world: Iteration 2 ($\gamma = 0.9$)

# Grid world: Iteration 3 ($\gamma = 0.9$)

# Grid world: Iteration 4 ($\gamma = 0.9$)

# Grid world: Iteration 5 ($\gamma = 0.9$)



| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| −0.4 | −0.19 | −0.15 | −0.16 | −0.15 | 0.54 | 0.91 | 1.55 | 0.87 | 0.3 |
| −0.19 | −0.04 | −0.03 | −0.03 | 0.64 | 0.94 | 1.77 | 2.23 | 1.74 | 0.87 |
| −0.15 | −0.03 | −0.02 | 0.41 | 0.74 | 1.65 | 2.24 | 3 | 2.23 | 1.55 |
| −0.16 | −0.04 | −0.11 | −0.53 | 0.57 | 0.95 | 1.79 | 2.24 | 2.18 | 0.91 |
| −0.18 | −0.09 | −0.54 | −5.48 | −0.53 | 0.64 | 0.96 | 2.62 | 2.7 | 2.09 |
| −0.16 | −0.06 | −0.14 | −0.66 | −0.13 | −0.01 | 2.47 | 3.32 | 5.51 | 3.72 |
| −0.17 | −0.06 | −0.24 | −1.16 | −0.22 | 2.23 | 3.21 | 5.97 | 7.52 | 6.08 |
| −0.21 | −0.15 | −1.07 | −10.96 | 0.52 | 2.39 | 5.5 | 7.47 | 10 | 7.8 |
| −0.24 | −0.12 | −0.29 | −1.14 | −0.25 | 2.2 | 3.19 | 5.94 | 7.54 | 6.07 |
| −0.4 | −0.24 | −0.28 | −0.41 | −0.26 | −0.18 | 2.16 | 3.38 | 5.49 | 3.88 |

# Grid world: Iteration 6 ($\gamma = 0.9$)

# Grid world: Iteration 7 ($\gamma = 0.9$)



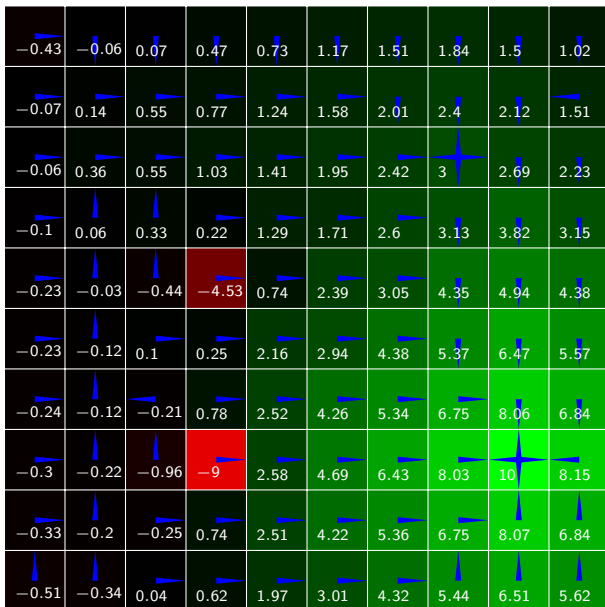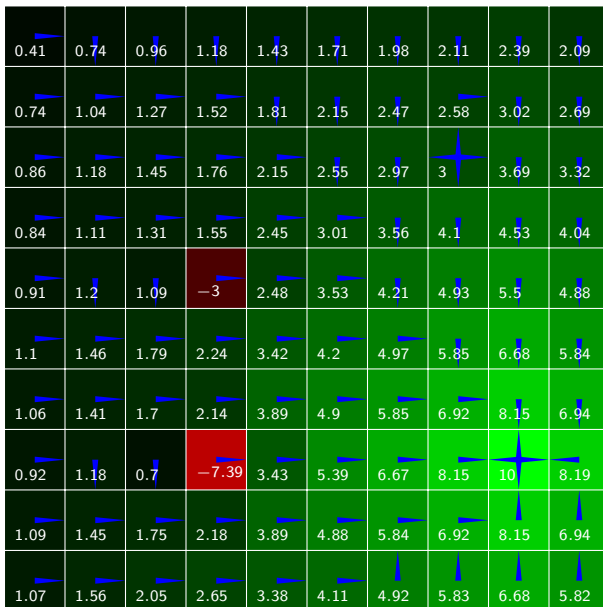| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| −0.41 | −0.21 | −0.17 | 0.17 | 0.37 | 0.96 | 1.34 | 1.75 | 1.31 | 0.78 |
| −0.21 | −0.06 | 0.27 | 0.42 | 1.04 | 1.38 | 1.94 | 2.35 | 1.94 | 1.31 |
| −0.17 | 0.13 | 0.24 | 0.8 | 1.18 | 1.84 | 2.36 | 3 | 2.39 | 1.81 |
| −0.18 | −0.06 | 0.09 | −0.11 | 0.96 | 1.38 | 2.09 | 2.48 | 3.17 | 2.13 |
| −0.21 | −0.11 | −0.58 | −5.15 | −0.05 | 1.6 | 2.09 | 3.75 | 4.22 | 3.66 |
| −0.2 | −0.09 | −0.18 | −0.7 | 1.19 | 1.8 | 3.74 | 4.74 | 6.19 | 5.1 |
| −0.21 | −0.1 | −0.28 | −0.14 | 1.32 | 3.58 | 4.7 | 6.52 | 7.92 | 6.65 |
| −0.27 | −0.19 | −1.05 | −10.02 | 1.8 | 3.9 | 6.15 | 7.88 | 10 | 8.07 |
| −0.29 | −0.18 | −0.34 | −0.14 | 1.28 | 3.52 | 4.7 | 6.5 | 7.94 | 6.65 |
| −0.46 | −0.3 | −0.36 | −0.46 | 0.86 | 1.77 | 3.59 | 4.85 | 6.23 | 5.21 |

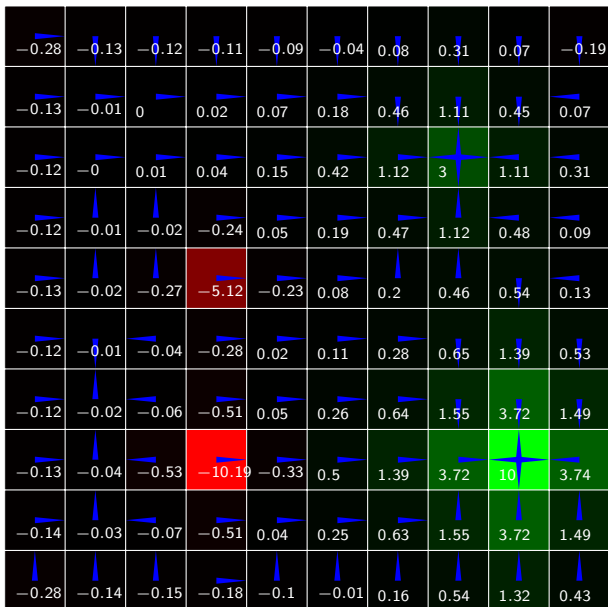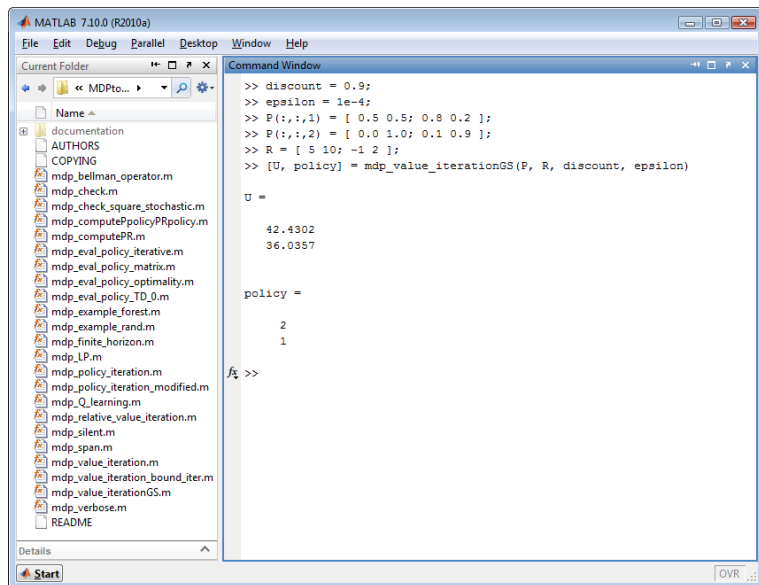# Grid world: Iteration 8 ($\gamma = 0.9$)

# Grid world: Iteration 9 ($\gamma = 0.9$)

# Grid world: Converged ($\gamma = 0.9$)

# Grid world: Converged ($\gamma = 0.5$)

# Markov Decision Process (MDP) Toolbox for Matlab
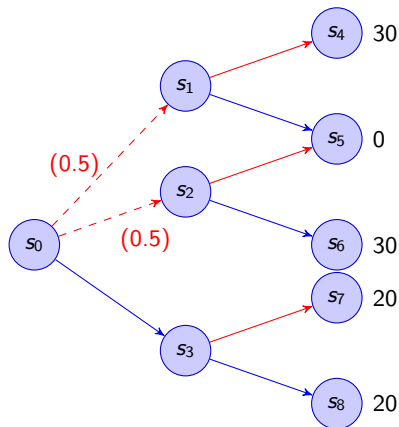
# Closed loop and open loop planning

- Closed loop: accounts for future state information (MDP)
- Open loop: does not account for future state information (path planning)



- Open loop plans do not always result in optimal behavior
- $U(r, r) = 15$
- $U(r, b) = 15$
- $U(b, r) = 20$
- $U(b, b) = 20$
- MDP solution can increase utility to 30

# Outline

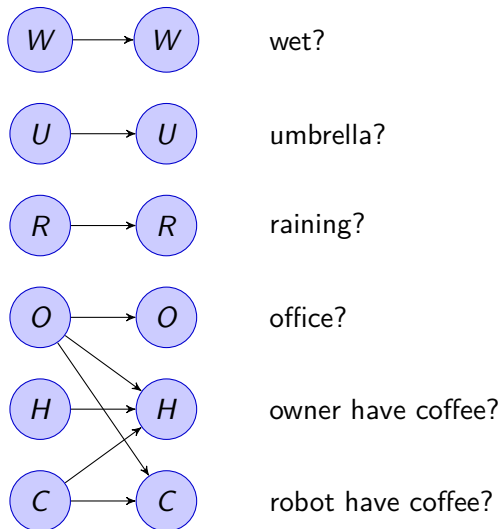# Factored Markov decision processes

- Number of discrete states grows exponentially with the number of state variables
- <span style="color:red">Factored Markov decision processes</span> (FMDPs) compactly represent transition and reward functions using dynamic decision networks
- Conditional probability distributions and utility functions can be represented as decision trees

# Coffee robot problem

- Robot must go buy a cup of coffee for its owner located at the office
- When it's raining, robot must get an umbrella to stay dry
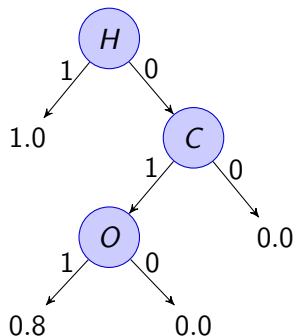- Actions: go to other location, buy coffee, deliver coffee, and get umbrella

# Transition model for deliver coffee action



W ——→ W     wet?

U ——→ U     umbrella?

R ——→ R     raining?

O ——→ O     office?

H ——→ H     owner have coffee?

C ——→ C     robot have coffee?

# Conditional distribution as decision tree

Representation of $P(H' = 1 \mid H, C, O, \text{deliver coffee})$

| H | C | O | $P(H')$ |
|---|---|---|---------|
| 1 | 1 | 1 | 1.0 |
| 1 | 1 | 0 | 1.0 |
| 1 | 0 | 1 | 1.0 |
| 1 | 0 | 0 | 1.0 |
| 0 | 1 | 1 | 0.8 |
| 0 | 1 | 0 | 0.0 |
| 0 | 0 | 1 | 0.0 |
| 0 | 0 | 0 | 0.0 |



Both the table and decision tree represent the same distribution, but the decision tree is more compact

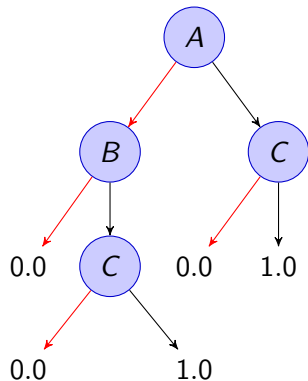# Structured dynamic programming

- $R(s, a)$ and $U(s)$ can also be represented using a decision tree
- Structured value iteration and structured policy iteration performs updates on leaves of the decision trees instead of all the states
- Structured dynamic programming algorithms improve efficiency by aggregating states
- Algorithms can also leverage additive decomposition of reward and value function

# Stochastic Planning Using Decision Diagrams (SPUDD)

Additional efficiency can be gained using decision diagrams



Software: `http://www.computing.dundee.ac.uk/staff/jessehoey/spudd/index.html`

DMU 4.3.2, J. Hoey, R. St-Aubin, A. Hu, *et al.*, "SPUDD: Stochastic planning using decision diagrams," in *Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence*, K. B. Laskey and H. Prade, Eds., Morgan Kaufmann, 1999, pp. 279–288

# Outline

# Approximate dynamic programming

- Discussion up to this point assumed (small) discrete state and action spaces
- Approximate dynamic programming is concerned with finding approximately optimal policies for problems with large or continuous spaces
- ADP is an active area of research and shares ideas with reinforcement learning
- Methods include parametric approximation, local approximation, and state aggregation (abstraction) methods

# Grid-based parametric approximation

- If state space $X$ is continuous, how to represent utility $U$?
- One way is to use a grid to get a discrete set of points $S \subset X$
- Let $V$ be a utility function over $S$
- Determine $U(x)$ by interpolating $V$
- Bellman update:
  $V_k(s) \leftarrow \max_a[R(s,a) + \gamma \sum_{x'} P(x' \mid s,a) U_{k-1}(x')]$



DMU 4.5.1, G. J. Gordon, "Approximate solutions to Markov decision processes," PhD thesis, Computer Science Department, Carnegie Mellon University, 1999

# Linear regression parametric approximation

- Instead of discretizing the state space (exponential in state variables) define a set of basis functions $\phi_1, \ldots, \phi_n$ over $X$
- Each basis function maps states in $X$ to real numbers
- Think of basis functions as a state-dependent features
- Represent $U$ as a linear combination of basis functions: $U(x) = \theta_1 \phi_1(x) + \ldots + \theta_n \phi_n(x)$
- Given fixed basis functions, determine $\theta_1, \ldots, \theta_n$ that best represents the optimal value function

# Linear regression parametric approximation

- Let $S$ be a finite set of samples from $X$ (doesn't have to be from a grid)
- Initialize $\theta_1, \ldots, \theta_n$ to 0, making $U_0(x) = 0$
- Bellman update:
  $V_k(s) \leftarrow \max_a [R(s, a) + \gamma \sum_{x'} P(x' \mid s, a) U_{k-1}(x')]$
- Use regression (e.g., least squares) to find $U_k$ that approximates $V_k$ at sample points in $S$
- Approach can be generalized to other parametric approximation methods (e.g., neural networks)

# Variable resolution dynamic programming

- ▶ Carve up $X$ into a discrete state space $S$
- ▶ Instead of solving for the optimal policy over $X$, solve for policy over $S$
- ▶ Some parts of the state space require finer resolution than other parts when finding a good policy
- ▶ One approach is to adaptively partition the state space according to heuristics

# Linear representations

Exact solutions for continuous state and action spaces are available under certain assumptions:

- Dynamics are linear Gaussian:
  $T(\mathbf{z} \mid \mathbf{s}, \mathbf{a}) = \mathcal{N}(\mathbf{z} \mid \mathbf{T}_s \mathbf{s} + \mathbf{T}_a \mathbf{a}, \Sigma)$
- Rewards are quadratic:
  $R(\mathbf{s}, \mathbf{a}) = \mathbf{s}^\top \mathbf{R}_s \mathbf{s} + \mathbf{a}^\top \mathbf{R}_a \mathbf{a}$
  assuming $\mathbf{R}_s = \mathbf{R}_s^\top \leq 0$ and $\mathbf{R}_a = \mathbf{R}_a^\top < 0$

Bellman equation (undiscounted) is:
$U_n(\mathbf{s}) = \max_{\mathbf{a}} \left( R(\mathbf{s}, \mathbf{a}) + \int T(\mathbf{z} \mid \mathbf{s}, \mathbf{a}) U_{n-1}(\mathbf{z}) d\mathbf{z} \right)$

# Linear representations

- It is possible to show that
  $U_n(\mathbf{s}) = \mathbf{s}^\top \mathbf{V}_n \mathbf{s} + q_n$
  $\mathbf{V}_n = \mathbf{T}_s^\top \mathbf{V}_{n-1} \mathbf{T}_s - \mathbf{T}_s^\top \mathbf{V}_{n-1} \mathbf{T}_a (\mathbf{T}_a^\top \mathbf{T}_a + \mathbf{R}_a)^{-1} \mathbf{T}_s^\top \mathbf{V}_{n-1} \mathbf{T}_s + \mathbf{R}_s$
  $q_n = q_{n-1} + \mathrm{Tr}(\Sigma \mathbf{V}_{n-1})$ assuming $\mathbf{R}_s = \mathbf{R}_s^\top \leq 0$ and
  $\mathbf{R}_a = \mathbf{R}_a^\top < 0$
- Optimal policy:
  $\pi_n(\mathbf{s}) = -(\mathbf{T}_a^\top \mathbf{V}_{n-1} \mathbf{T}_a + \mathbf{R}_a)^{-1} \mathbf{T}_a \mathbf{V}_{n-1} \mathbf{T}_s \mathbf{s}$
- Note that optimal policy does not depend upon noise!
- For some problems, the dynamics are approximated with a linear Gaussian system and then solved exactly

# Direct policy search

- Alternative to planning over entire state space is to search the space of policies directly
- Consider stochastic policy parameterized by $\lambda$
- Probability policy selects $a$ in state $s$ is given by $\pi_\lambda(a \mid s)$
- Use Monte Carlo simulations to estimate expected discounted return from a distribution of initial states
- Find best $\lambda$ using local search, cross entropy method, or evolutionary algorithms

# Outline

1. Formulation
2. Dynamic programming
3. Structured dynamic programming
4. Approximate dynamic programming
5. Online methods

# Online methods

- Online methods compute optimal action from current state
- Expand tree up to some horizon
- States reachable from the current state is typically small compared to full state space
- Heuristics and branch-and-bound techniques allow search space to be pruned
- Monte Carlo methods provide approximate solutions

## Forward search

Provides optimal action from current state $s$ up to depth $d$

Recall $U(s) = \max_{a \in A(s)}[R(s, a) + \gamma \sum_{s' \in S(s,a)} T(s' \mid s, a)U(s')]$

```
1: function SELECTACTION(s, d)
2:     if d = 0
3:         return (NIL, 0)
4:     (a*, v*) ← (NIL, -∞)
5:     for a ∈ A(s)
6:         q ← R(s, a)
7:         for s' ∈ S(s, a)
8:             (a', v') ← SELECTACTION(s', d - 1)
9:             q ← q + γT(s' | s, a)v'
10:        if q > v*
11:            (a*, v*) ← (a, q)
12:    return (a*, v*)
```

Time complexity is $O((|S| \times |A|)^d)$

# Branch and bound search

Requires a lower bound $\underline{U}(s)$ and upper bound $\overline{U}(s, a)$

```
 1: function SELECTACTION(s, d)
 2:     if d = 0
 3:         return (NIL, U(s))
 4:     (a*, v*) ← (NIL, −∞)
 5:     for a ∈ A(s)
 6:         if U(s, a) < v*
 7:             return (a*, v*)
 8:         q ← R(s, a)
 9:         for s' ∈ S(s, a)
10:             (a', v') ← SELECTACTION(s', d − 1)
11:             q ← q + γT(s' | s, a)v'
12:         if q > v*
13:             (a*, v*) ← (a, q)
14:     return (a*, v*)
```
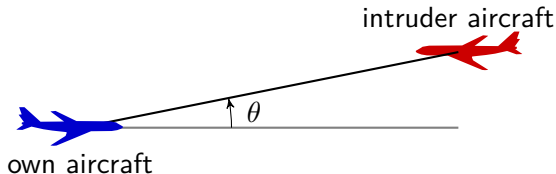
# Sparse sampling

Requires a generative model $(s', r) \sim G(s, a)$

```
 1: function SELECTACTION(s, d)
 2:     if d = 0
 3:         return (NIL, 0)
 4:     (a*, v*) ← (NIL, −∞)
 5:     for a ∈ A(s)
 6:         q ← 0
 7:         for i ← 1 to n
 8:             (s', r) ∼ G(s, a)
 9:             (a', v') ← SELECTACTION(s', d − 1)
10:             q ← q + (r + γv')/n
11:         if q > v*
12:             (a*, v*) ← (a, q)
13:     return (a*, v*)
```
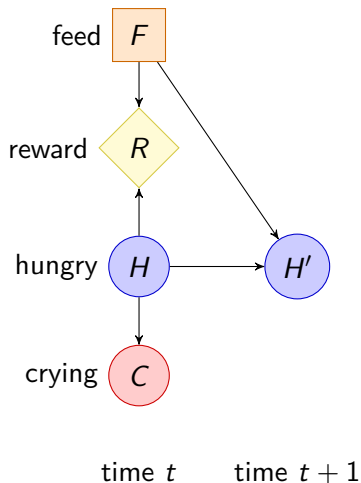
# Example POMDP

Aircraft collision avoidance using an electro-optical sensor



- ▶ Own aircraft can accelerate up or down
- ▶ Large penalty for collision, small penalty for maneuvering
- ▶ Dynamics are a function of current position and velocity
- ▶ Electro-optical sensor only measures angle
- ▶ Limited field of view

Optimization resulted in "passive ranging" behavior.

S. Temizer, M. J. Kochenderfer, L. P. Kaelbling, T. Lozano-Pérez, and J. K. Kuchar, "Collision Avoidance for Unmanned Aircraft using Markov Decision Processes," in AIAA Guidance, Navigation, and Control Conference, Toronto, Canada, 2010.

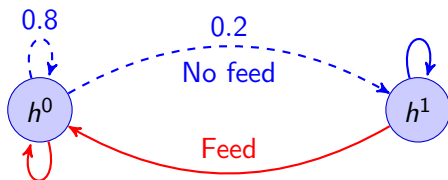# Crying baby problem



- Need to decide whether to feed baby given whether baby is crying
- Crying is a noisy indication that the baby is hungry
- Costs 5 to feed and 10 for hungry
- Infinite horizon with discount $\gamma$

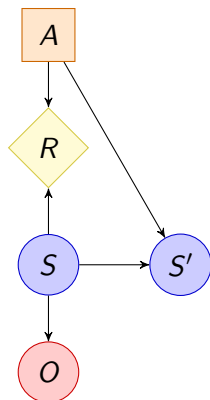# Crying baby problem

Transition model



- $P(c^1 \mid h^0) = 0.1$ (cry when not hungry)
- $P(c^1 \mid h^1) = 0.8$ (cry when hungry)

# POMDPs



- POMDP = MDP + sensor model
- Sensor model: $O(o \mid s)$ or sometimes $O(o \mid s, a)$
- Decisions can only be based on history of observations $o_1, o_2, \ldots, o_t$
- Instead of keeping track of arbitrarily long histories, we keep track of the belief state
- A belief state is a distribution over states; in belief state $b$, probability $b(s)$ is assigned to being in $s$

Policies in POMDPs are mappings from belief states to actions

# Computing belief states

- Begin with some initial belief state $b$ prior to any observations
- Compute new belief state $b'$ based on current belief state $b$, action $a$, and observation $o$
- $b'(s') = P(s' \mid o, a, b)$
- $b'(s') \propto P(o \mid s', a, b)P(s' \mid a, b)$
- $b'(s') \propto O(o \mid s', a)P(s' \mid a, b)$
- $b'(s') \propto O(o \mid s', a) \sum_s P(s' \mid a, b, s)P(s \mid a, b)$
- $b'(s') \propto O(o \mid s', a) \sum_s T(s' \mid s, a)b(s)$
- Kalman filter: exact update of the belief state for linear dynamical systems
- Particle filter: approximate update for general systems
- For now, we're only considering discrete state problems

# Crying baby example

- $b = (h^0, h^1) = (0.5, 0.5)$
- No feed, cry
- $b = (0.0928, 0.9072)$
- Feed, no cry
- $b = (1, 0)$
- No feed, no cry
- $b = (0.9759, 0.0241)$
- No feed, no cry
- $b = (0.9701, 0.0299)$
- No feed, cry
- $b = (0.4624, 0.5376)$

# POMDP execution

1. Initialize belief state $b$
2. Execute $a = \pi(b)$
3. Observe $o$
4. Update $b$ based on $b$, $a$, and $o$
5. Go to 2

# POMDP as belief MDP

- POMDP is really an MDP where the states are belief states
- $\mathcal{B}$, the set of belief states, comprise the state space
  - If there are $n$ discrete states, then $\mathcal{B} \subset \mathbb{R}^n$
- $\mathcal{A}$, the set of actions, remains the same
- $\tau(b' \mid b, a)$ is the state transition function
  - See DMU equation 6.5
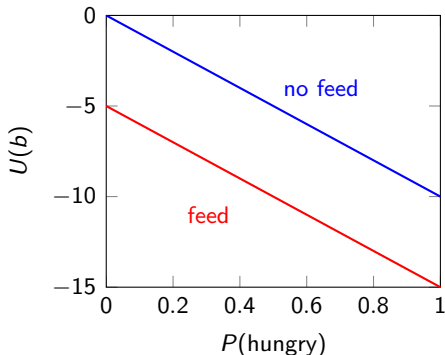- $\rho(b, a)$, the immediate reward, is equal to $\sum_s b(s) R(s, a)$

Discrete MDP methods cannot be used directly for POMDPs

# Outline

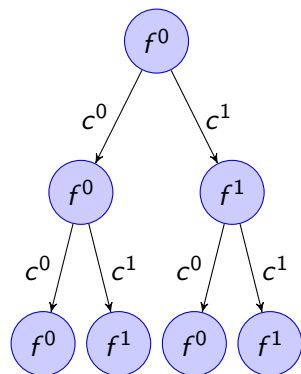# Alpha vectors

- For now, assume a 1-step horizon
- $U(s) = \max_a R(s, a)$
- $U(b) = \max_a \sum_s b(s) R(s, a)$
- Let $\alpha_a$ be $R(\cdot, a)$ represented as a vector
- $U(b) = \max_a b \cdot \alpha_a$
- The alpha vectors define hyperplanes in belief space
- One-step utility is <span style="color:red">piecewise linear and convex</span>
- What about $H$-step utility?

# Conditional plans
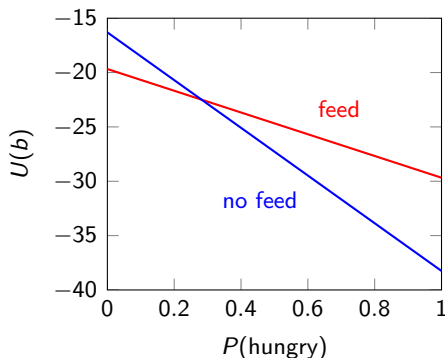
### Example 3-step conditional plan



- A conditional plan specifies what to do from the initial belief state after each possible observation up to some horizon
- Let $\alpha_p$ represent the expected utility of following plan $p$
- $U(b) = \max_p b \cdot \alpha_p$
- $H$-step utility is also piecewise linear and convex
- Execute action at root node of plan that maximizes $b \cdot \alpha_p$

# POMDP value iteration

- Value iteration for POMDPs involves generating conditional plans and associated alpha vectors up to some depth

- Number of possible conditional plans is doubly exponential

- For the crying baby problem, there are $2^{255}$ conditional 8-step plans

- During value iteration, toss out dominated plans



Optimal crying baby policy

Exact solution for general POMDP is PSPACE-hard

# Outline

# Offline methods

- Offline POMDP methods involve doing all or most of the computing prior to execution
- Practical method generally find only approximate solutions
- Some methods output policy in terms of alpha-vectors, others as finite-state controllers
- Some methods involve leveraging a factored representation
- Very active area of research

# QMDP value approximation

- QMDP approximation method involves solving for the MDP assuming full observability
- Compute $Q(s, a)$ for all states and actions
- Let $\alpha_a$ be the vector representation of the $Q$ function for action $a$
- QMDP utility function: $U_{\mathsf{QMDP}}(b) = \max_a b \cdot \alpha_a$
- QMDP policy: $\pi_{\mathsf{QMDP}}(b) = \arg\max_a b \cdot \alpha_a$

QMDP assumes all state uncertainty disappears at the next step

# QMDP value approximation

- Upper bound on utility:
  $U_{\mathsf{QMDP}}(b) \geq U^*(b)$
- Very small change in policy for crying baby problem
- QMDP tends to have difficulty with problems with information gathering actions (e.g., look over right shoulder when changing lanes)
- Performs extremely well in many real problems (e.g., airborne collision avoidance)



Crying baby policy

Dashed lines are for QMDP

# Fast informed bound (FIB) value approximation

- FIB computes single alpha vector for each action (just like QMDP)
- Tighter bounds are provided by taking into account observation model
- QMDP:
  $\alpha_a^{(k+1)}(s) = R(s, a) + \gamma \sum_{s'} T(s' \mid s, a) \max_{a'} \alpha_{a'}^{(k)}(s')$
- FIB:
  $\alpha_a^{(k+1)}(s) = R(s, a) + \gamma \sum_o \max_{a'} \sum_{s'} O(o \mid s', a) T(s' \mid s, a) \alpha_{a'}^{(k)}(s')$
- FIB complexity is $O(|A|^2 |S|^2 |O|)$
- FIB is also an upper bound on optimal value function

# Point-based value iteration methods

- Point-based value iteration involves backing up alpha vectors $\Gamma = \boldsymbol{\alpha}_{1:n}$ associated with belief points $\mathbf{b}_{1:n}$
- Choose points $\mathbf{b}_{1:n}$ and initialize alpha vectors to lower bound
- Update alpha vectors in $\Gamma$ at those points such that $U^\Gamma(\mathbf{b}) \leq U^*(\mathbf{b})$

```
1: function BACKUPBELIEF(Γ, b)
2:     for a ∈ A
3:         for o ∈ O
4:             b' ← UPDATEBELIEF(b, a, o)
5:             α_{a,o} ← arg max_{α∈Γ} α^⊤ b'
6:         α_a(s) ← R(s, a) + γ ∑_{s',o} O(o | s', a) T(s' | s, a) α_{a,o}(s')
7:     α ← arg max_{α_a} α_a^⊤ b
8:     return α
```

# Point-based value iteration methods

- Some algorithms use a fixed grid of belief states
- Better to focus effort on reachable belief states under random policies
- Even better to focus effort on belief states reachable under the optimal policy
- Successive Approximations of the Reachable Space under Optimal Policies (SARSOP) algorithm is one of the best approximate solvers for POMDPs

DMU 6.4.3, MDPAI 7.5, 7.6, AIMA 17.7, H. Kurniawati, D. Hsu, and W. Lee, "SARSOP: efficient point-based POMDP planning by approximating optimally reachable belief spaces," in *Robotics: Science and Systems*, 2008

# APPL: Approximate POMDP Planning Toolkit

- APPL contains an efficient C++ implementation of SARSOP
- Supports factored (and non-factored) representations
- Takes as input a POMDP model (in text or XML) and outputs a policy represented as XML
- Source is available (g++ and Visual Studio)
- Includes tools for simulation, evaluation, and visualization

# Example input file

```
discount: 0.9
values: reward
states: h0 h1
actions: f0 f1
observations: c0 c1
start: uniform

T: f0 : h0 : h1 0.1
T: f0 : h0 : h0 0.9
T: f0 : h1 : h1 1
T: f1 : * : h0 1

O: * : h0
0.9 0.1

O: * : h1
0.2 0.8

R: f0 : h1 : * : * -10
R: f1 : h0 : * : *  -5
R: f1 : h1 : * : * -15
```

## Example execution

```
Loading the model ...
  input file   : cry.pomdp
  loading time : 0.00s

SARSOP initializing ...
  initialization time : 0.00s

--------------------------------------------------------------------------------
 Time   |#Trial |#Backup |LBound   |UBound    |Precision  |#Alphas |#Beliefs
--------------------------------------------------------------------------------
 0.001   0       0        -55.0001  -22.7678   32.2323     2        1
 0.002   5       51       -25.3301  -24.5726   0.757554    9        6
 0.003   10      105      -24.6886  -24.6415   0.0471187   2        11
 0.004   13      150      -24.6758  -24.6552   0.0205842   5        11
 0.005   17      207      -24.675   -24.6699   0.0050443   2        16
 0.006   20      259      -24.6749  -24.6726   0.00228751  2        18
 0.007   23      305      -24.6749  -24.674    0.000950473 2        18
 0.007   23      305      -24.6749  -24.674    0.000950473 2        18
--------------------------------------------------------------------------------

Writing out policy ...
  output file : out.policy
```

# Example output file

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<Policy version="0.1" type="value" model="cry.pomdp"
xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
xsi:noNamespaceSchemaLocation="policyx.xsd">
<AlphaVector vectorLength="2" numObsValue="1" numVectors="2">
<Vector action="0" obsValue="0">-16.3055 -38.2512 </Vector>
<Vector action="1" obsValue="0">-19.6749 -29.6749 </Vector>
</AlphaVector> </Policy>
```

# Outline

# Online methods

- Online methods determine the optimal policy by planning from the current belief state
- Many online methods use a depth-first tree-based search up to some horizon
- The belief states reachable from the current state is typically small compared to the full belief space
- Time complexity is exponential in the search depth
- Heuristics and brand-and-bound techniques allow search space to be pruned

## Forward search

```
 1: function SELECTACTION(b, d)
 2:     if d = 0
 3:         return (NIL, U(b))
 4:     (a*, u*) ← (NIL, -∞)
 5:     for a ∈ A
 6:         u ← R(b, a)
 7:         for o ∈ O
 8:             b' ← UPDATEBELIEF(b, a, o)
 9:             (a', u') ← SELECTACTION(b', d - 1)
10:             u ← u + γP(o | b, a)u'
11:         if u > u*
12:             (a*, u*) ← (a, u)
13:     return (a*, u*)
```

## Branch and bound

```
1: function SELECTACTION(b, d)
2:     if d = 0
3:         return (NIL, U(b))
4:     (a*, u) ← (NIL, -∞)
5:     for a ∈ A
6:         if U̅(b, a) ≤ u
7:             return (a*, u)
8:         u ← R(b, a)
9:         for o ∈ O
10:            b' ← UPDATEBELIEF(b, a, o)
11:            (a', u') ← SELECTACTION(b', d - 1)
12:            u ← u + γP(o | b, a)u'
13:        if u > u
14:            (a*, u) ← (a, u)
15:    return (a*, u)
```

# Online methods

- In practice, online methods can be very slow, requiring relatively shallow horizons
- Instead of single-step actions, can use action sequences or control strategies to reduce required depth
- Hybrid methods combine offline and online methods; offline approximate solution provides rough heuristics for leaf nodes in online search
- Online methods are usually insensitive to dimensionality of the state space

DMU 6.5, S. Ross, J. Pineau, S. Paquet, *et al.*, "Online planning algorithms for POMDPs," *Journal of Artificial Intelligence Research*, vol. 32, pp. 663–704, 2008

# Some things we didn't talk about

- Partially observable stochastic games (POSGs)
  - See MDPAI 8
- Decentralized POMDPs
  - See MDPAI 9
- Reinforcement learning in POMDPs
  - See S. Thrun, "Monte Carlo POMDPs," in *Advances in Neural Information Processing Systems*, S. A. Solla, T. K. Leen, and K.-R. Müller, Eds., vol. 12, Cambridge, Mass.: MIT Press, 2000, pp. 1064–1070
- Continuous state POMDPs
  - See H. Bai, D. Hsu, W. S. Lee, *et al.*, "Monte Carlo value iteration for continuous state POMDPs," in *International Workshop on the Algorithmic Foundations of Robotics*, 2010
- Mostly observable MDPs
  - See S. C. W. Ong, S. W. Png, D. Hsu, *et al.*, "Planning under uncertainty for robotic tasks with mixed observability," *International Journal of Robotics Research*, vol. 29, no. 8, pp. 1053–1068, 2010