

# Partially Observable Markov Decision Processes POMDP

Marius Bulacu

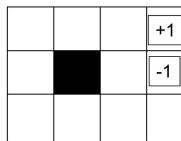


Kunstmatige Intelligentie / RuG



## Example

- Agent doesn't know where it is
  - Has no sensors
  - Being put in an unknown square
- What should the agent do?
  - If it knew that it is in (3,3) it would do *Right* action
  - But it doesn't know

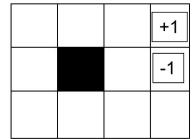


3

4

## Solutions

- First solution:
  - Step 1: reduces uncertainty
  - Step 2: try heading to the +1 exit
- Reduce uncertainty
  - move 5 times *Left* so it is quite likely to be at the left wall
  - Then move 5 times *Up* so it is quite likely to be at left top wall
  - Then move 5 time *Right* to goal state
  - Continue moving right to increase chance to get to +1
- Quality of solution
  - Chance of 81.8% to get to +1
  - Expected utility of about 0.08
- Second solution: POMDP



5

6

## MDP

- Components:
  - States - S
  - Actions - A
  - Transitions -  $T(s, a, s')$
  - Rewards –  $R(s)$
- Problem:
  - choose the action that makes the right tradeoffs between the immediate rewards and the future gains, to yield the best possible solution
- Solution:
  - Policy

## PO-MDP

- Components:
  - States - S
  - Actions - A
  - Transitions –  $T(s,a,s')$
  - Rewards –  $R(s)$
  - Observations –  $O(s,o)$

7

8

## Policy Mapping: MDP vs POMDP

- In MDP, mapping is from states to actions.
- In POMDP, mapping is from probability distributions (over states) to actions.

## Observation Model

- $O(s, o)$  – probability of getting observation  $o$  in state  $s$
- Example:
  - for robot without sensors:
    - $O(s, o) = 1$  for every  $s$  in  $S$
    - Only one observation  $o$  exists

## Belief State

- The agent may be in any state
- $b(s)$  = probability distribution over all states
- Example:
  - in the 4x3 example for robot without sensors the initial state:

$$b = \langle 1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 1/9, 0, 0 \rangle$$

## Belief State (contd.)

- $b(s)$  = probability assigned to actual state  $s$  by belief  $b$ 
  - $0 < b(s) < 1$  for every  $s$  in  $S$
  - sum of  $b(s)$  over all states = 1
- We will define  $b' = \text{FORWARD}(b, a, o)$

$$b'(s') = \alpha O(s', o) \sum_s T(s, a, s') b(s)$$

$\alpha$  is normalizing constant that makes the belief state sum to 1

## Finding the Best Action

- The optimal action depend only on the current belief state
- Choosing the action:
  - Given belief state  $b$ , execute  $\pi^*(b)$
  - Received observation  $o$
  - Calculate new state belief  $b' = \text{FORWARD}(b, a, o)$
  - Set  $b \leftarrow b'$

## $\pi^*(b)$ vs $\pi^*(s)$

- POMDP belief state is continuous
- In the 4x3 example,  $b$  is in an 11-dimensional continuous space
  - $b(s)$  is a point on a line between 0 to 1
  - $b$  is point in  $n$  dimensional space

## Transition Model & Rewards

- Need different view of transition model and reward
- Transition model as  $T(b, a, b')$ 
  - Instead  $T(s, a, s')$
  - Actual state  $s$  is not known
- Reward as  $R(b)$ 
  - Instead  $R(s)$
  - Have to model the uncertainty in the belief state

## Transition Model

- Probability of getting the observation  $o$ 

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

$$= \sum_{s'} O(s', o) P(s' | a, b) = \sum_{s'} O(s', o) \sum_s T(s, a, s') b(s)$$
- Transition model
 
$$t(b, a, b') = P(b' | a, b) = \sum_o P(b' | o, a, b) P(o | a, b)$$

$$= \sum_o P(b' | o, a, b) \sum_{s'} O(s', o) \sum_s T(s, a, s') b(s)$$

where  $P(b' | o, a, b)$  is 1 if  $b' = \text{FORWARD}(b, a, o)$  and 0 otherwise

## Reward Function

## From POMDP to MDP

- Expected reward of all the states the agent might be in

$$R(b) = \sum_s b(s) R(s)$$

- $T(b, a, b')$  and  $R(b)$  define an MDP
  - $\pi^*(b)$  for this MDP is optimal policy for original POMDP
  - Belief state is observable to the agent
- Need new versions of Value / Policy Iteration
  - for the continuous belief state

## Back to the Example

---

- POMDP solution for the 4x3 environment:  
[Left, Up, Up, Right, Up, Up, Right, Up, Up, Right, Up, Right, Up, Right, Up ...]
- The policy is a sequence since the problem is deterministic in beliefs space
- The agent gets to the goal 86.6% of times
  - Expected utility is 0.38

## Overview of Markov Models

---

Markov Models	Do we have control over the state transitions?	
	NO	YES
Are the states completely observable?	YES	<b>Markov Chain</b>
	NO	<b>HMM</b> Hidden Markov Model
		<b>MDP</b> Markov Decision Process
		<b>POMDP</b> Partially Observable Markov Decision Process