

Partially observable Markov decision processes

Jiří Kléma

Department of Computer Science,
Czech Technical University in Prague



<https://cw.fel.cvut.cz/wiki/courses/b4b36zui/prednasky>

Agenda

- Previous lecture: Markov decision processes (MDPs)
 - stochastic process with a limited memory,
 - world/environment well defined by its transition and reward functions,
 - goal to find the optimal policy,
 - dynamic programming most frequently used,
- partially observable Markov decision processes (POMDPs)
 - the world is partially observable only, states are not available,
 - define a new stochastic process that generalizes MDP,
 - policy changes, complexity grows, theoretical and real solutions.

Partial observability

- MDPs work with the assumption of complete observability
 - assumption that the actual state s is always known is often non realistic,
 - examples
 - * physical processes such as a nuclear reactor, complex machines,
 - * we do know the physical laws that underlie the process,
 - * we know the structure and characteristics of the machine and its parts,
 - * however, do not know the initial state and subsequent states, can only measure temperature,
 - * or have signals from various (unreliable) sensors.

Partial observability

- partially observable Markov decision process (POMDP)
 - MDP generalization, states guessed from observations coupled with them,
 - POMDP = $\{S, A, P, R, O, \Omega\}$,
 - * O is a set of observations,
 - * Ω is a sensoric model that defines conditional observation probs

$$\Omega_{s'o}^a = Pr\{o_{t+1} = o \mid s_{t+1} = s', a_t = a\}$$

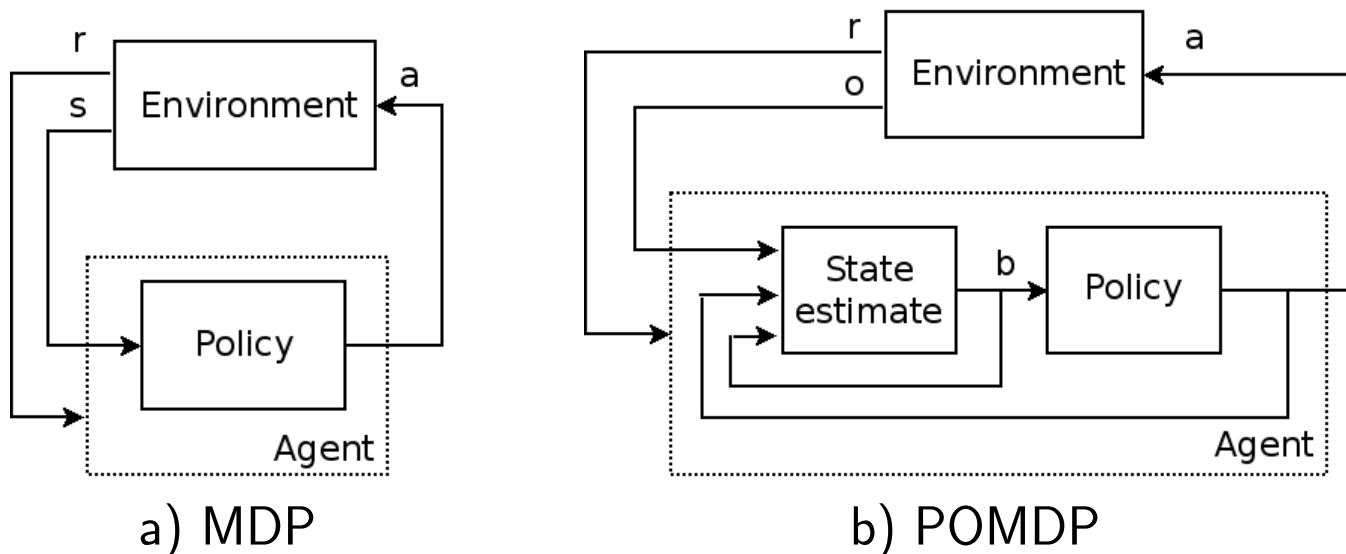
- instead of s agent internally keeps prob distribution b (**belief**) across states
 - * we perform a in unknown s (knowing $b(s)$ only) and observe o ,
 - * then we update our belief

$$b'(s') = \eta \Omega_{s'o}^a \sum_{s \in S} P_{ss'}^a b(s)$$

- * η is a normalization constant such that $\sum_{s' \in S} b'(s') = 1$.

Partial observability

- consequences of partial observability
 - it makes no sense to concern policy $\pi : S \rightarrow A$, shift to $\pi : B \rightarrow A$,
 - commonly computationally intractable, approximate solutions only
 - * for n states, b is an n -dimensional real vector,
 - * PSPACE-hard, worse than NP.



Partial observability – example

- b space is 1D $\rightarrow V(b)$ is a real function of one variable,
- assumed that in near points of b space will be
 - very similar utility and identical policy,
- policy is equivalent to a **conditional plan** dependent on future observations
 - example, plan of length 2: $[Stay, \text{if } O = o_0 \text{ then } Go \text{ else } Stay]$,
- let $\alpha_p(s)$ be the utility of plan p starting from state s
 - then the same plan executed from b has the utility

$$\sum_s b(s)\alpha_p(s) = b \cdot \alpha_p$$

- α_p is a linear function of b (hyperplane for complex spaces),
- optimal policy follows the plan with highest expected utility

$$V(b) = V^{\pi^*}(b) = \max_p b \cdot \alpha_p$$

- $V(b)$ is a partially linear function of b .

Partial observability – example

:: there are two plans of length 1, for them it holds

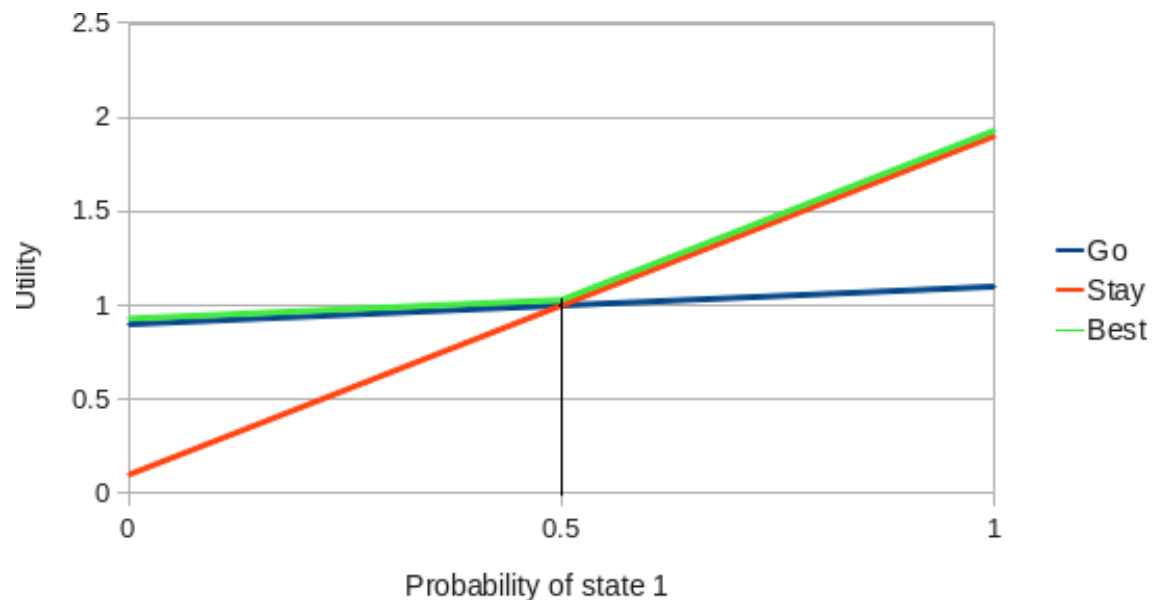
$$\alpha_{[Stay]}(0) = R(0) + \gamma(.9R(0) + .1R(1)) = 0.1$$

$$\alpha_{[Stay]}(1) = R(1) + \gamma(.9R(1) + .1R(0)) = 1.9$$

$$\alpha_{[Go]}(0) = R(0) + \gamma(.9R(1) + .1R(0)) = 0.9$$

$$\alpha_{[Go]}(1) = R(1) + \gamma(.9R(0) + .1R(1)) = 1.1$$

$$\alpha_{[Stay]}(b(1) = 0.3) = .7\alpha_{[Stay]}(0) + .3\alpha_{[Stay]}(1) = 0.64$$



Partial observability – example

:: there are 8 plans of length 2 (4 dotted plans dominated by other plans)

[*Stay*, if $O = o_0$ then *Go* else *Stay*] encoded as [SGS]

$$\begin{aligned}
 \alpha_{[SSS]}(0) &= R(0) + \gamma(.9\alpha_{[S]}(0) + .1\alpha_{[S]}(1)) &= 0.28 \\
 \alpha_{[SSS]}(1) &= R(1) + \gamma(.9\alpha_{[S]}(1) + .1\alpha_{[S]}(0)) &= 2.72 \\
 \alpha_{[SGS]}(0) &= R(0) + \gamma(.9(.6\alpha_{[G]}(0) + .4\alpha_{[S]}(0)) + .1(0.4\alpha_{[G]}(1) + .6\alpha_{[S]}(1)) &= 0.68 \\
 \alpha_{[SGS]}(1) &= R(1) + \gamma(.9(.4\alpha_{[G]}(1) + .6\alpha_{[S]}(1)) + .1(0.6\alpha_{[G]}(0) + .4\alpha_{[S]}(0)) &= 2.48
 \end{aligned}$$

