

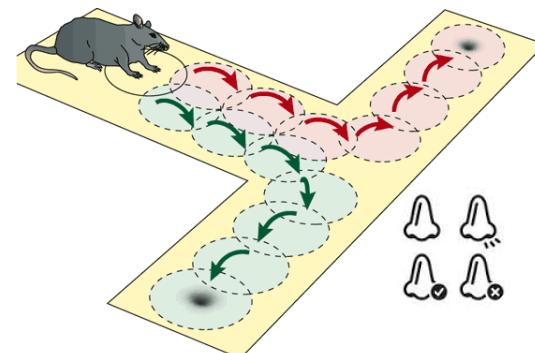
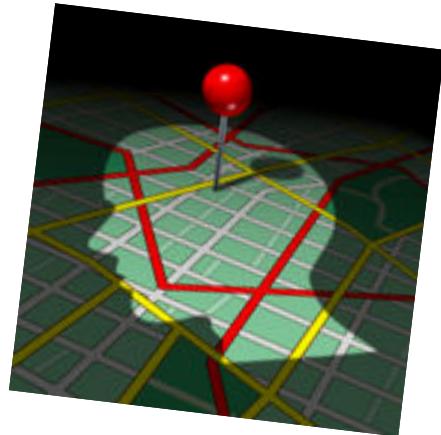


Partially Observable Markov Decision Processes (POMDP)

Lionel Rigoux & Frederike Petzschner

Introduction

- MDP >> Full observability: the agent always knows the state of the world
- This might often not be true in real life
 - *Imperfect memory*
// navigation: “turn left on the seventh street”
> what if you loose track of the number of streets already passed?
 - *Changing environment*
// reward selection in a T-maze
> smell can predict outcomes



Outline

- Extend the MPD framework to account for state uncertainty
 - Beliefs representation
 - Observation function
 - Belief updating and state chaining
- Formalization
- Solution
- Conclusion
- Perspectives



state**action**

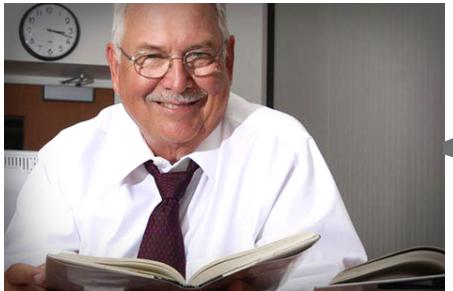
leave

stay

stay

stay

leave

reward $R = 100$  $R = 30$ 

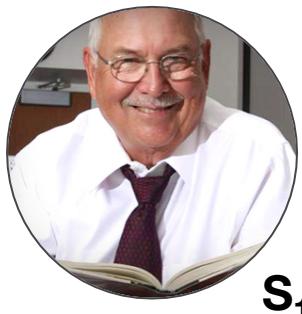
$$\begin{aligned} v_{fx}^s &= f(x^s) + \gamma \text{value}(y^s) \\ \nabla_x(f^s) &= (x_f)^s + f(\nabla_x s) \\ s &= f^{-1}\pi_i \rightarrow \nabla_x s = (x_f)\pi_i \\ \nabla_x \pi_i &= A^{\pi_k} \cdot \pi_k \\ &\rightarrow (x_f + A^{\pi_k} \cdot \pi_k) \end{aligned}$$

 $R = -40$ 

leave

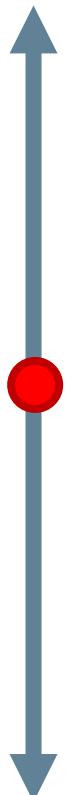


state
not known



belief
 $b=p(s=S_1)$

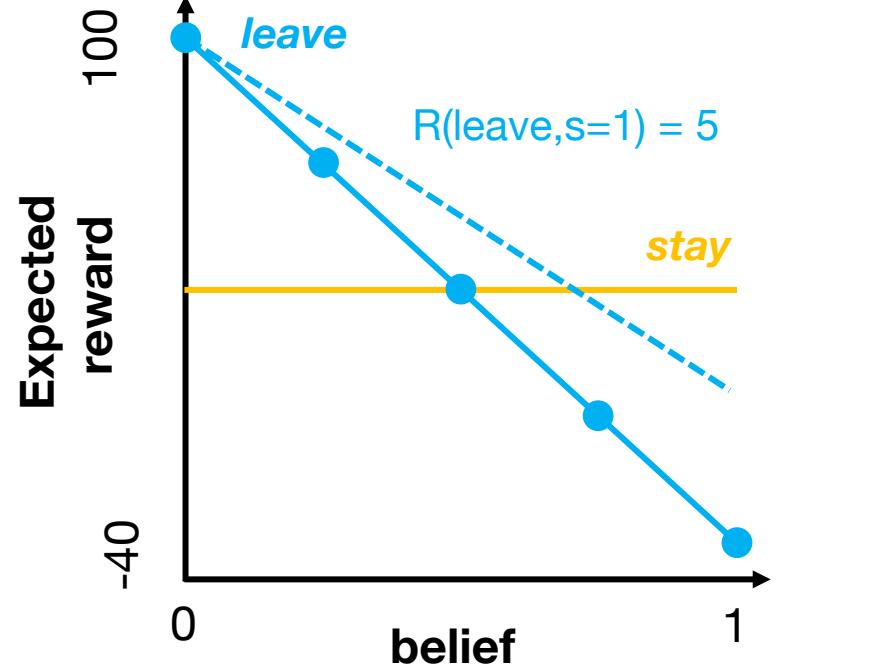
$$p(s=S_1) = 0$$



$$p(s=S_1) = 1$$

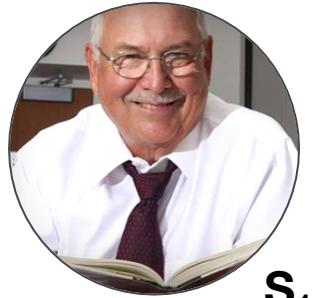
actions and payoff function

optimal policy



$$E[R](a) = p(s=0) R_0(a) + p(s=1) R_1(a)$$





observation function

provide information about state

	<i>leave</i>	<i>stay</i>	<i>listen</i>
<i>noises</i>	0	0.5	0.15
<i>no one</i>	1	0.5	0.85

$$b' = \frac{p(o|s', a) \sum_s p(s'|s, a) b(s)}{\sum_{s'} p(o|s', a) \sum_s p(s'|s, a) b(s)}$$

	<i>leave</i>	<i>stay</i>	<i>listen</i>
<i>noises</i>	1	0.5	0.85
<i>no one</i>	0	0.5	0.15



$$p(s=S_1) = 0$$

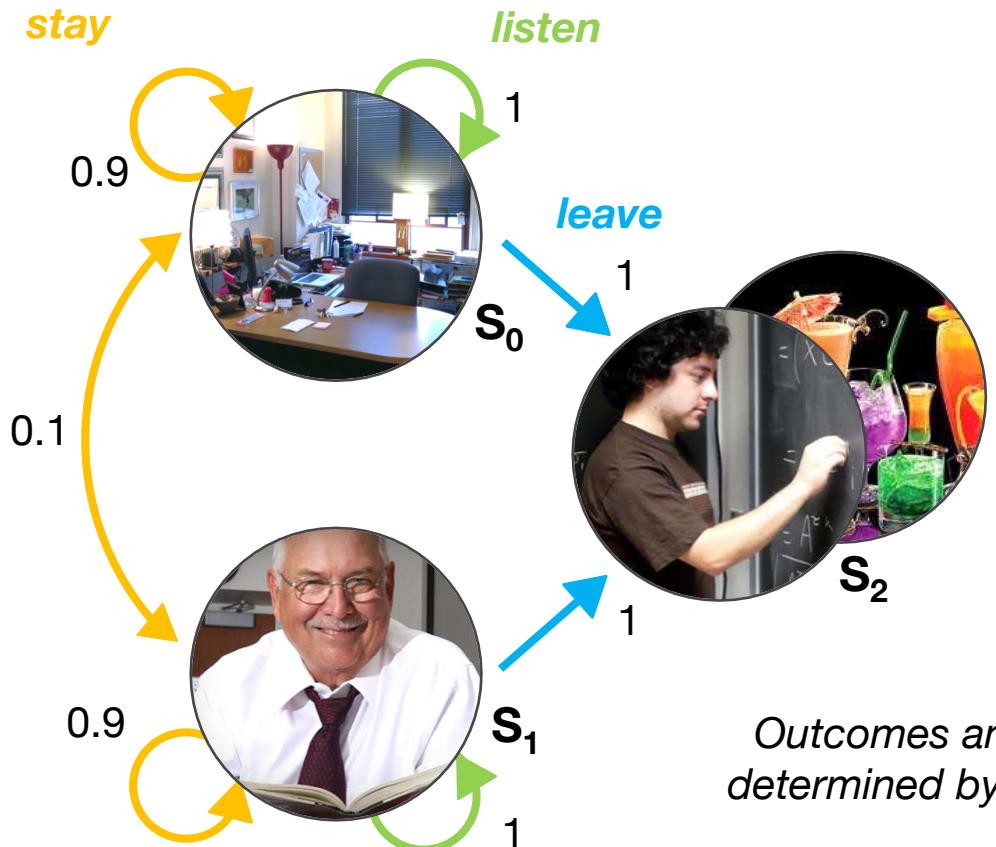


listen
no one

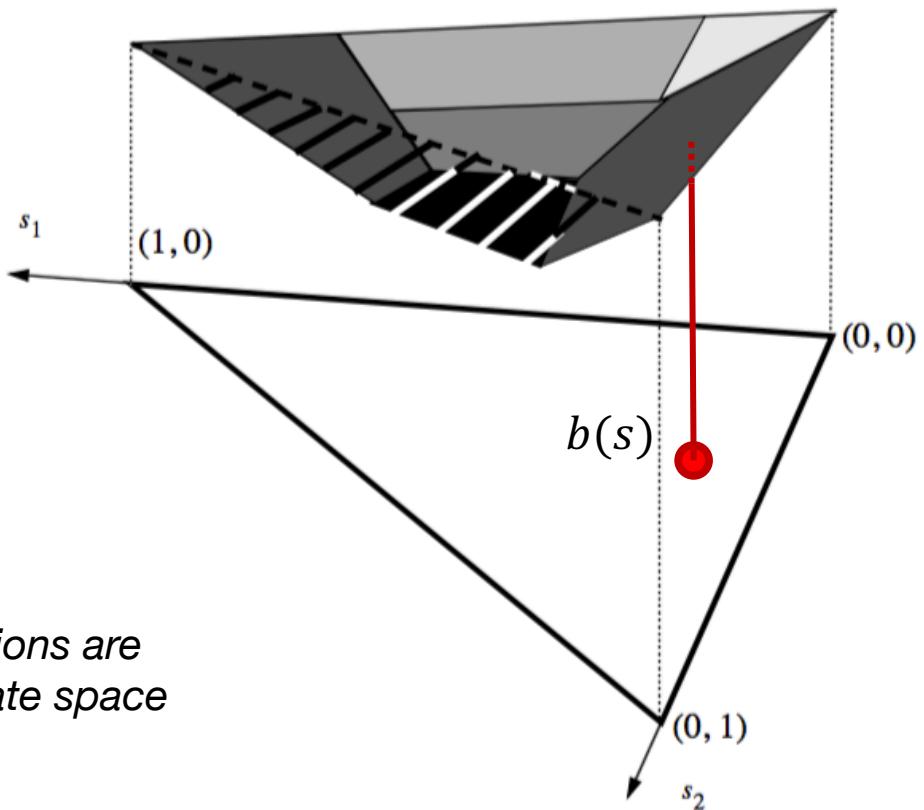
$$p(s=S_1) = 1$$



state space



belief space



Outcomes and observations are determined by the real state space

Policy relies on the belief state



POMDP Formalism

MDP

- S set of states
- A set of actions
- T transition matrix $S \times A \rightarrow S$
- R reward function $S \times A \rightarrow \mathbb{R}$
- γ discount factor

POMDP extension

- Ω set of observations
- O observation probabilities
 $S \times A \times \Omega \rightarrow [0, 1]$
- B belief space
- r reward function $B \times A \rightarrow \mathbb{R}$
- τ belief update function $B \times A \times \Omega \rightarrow B$

Simulation workflow

- Initial state (s, b)
- Select action $a = \pi(b)$
- Update state $s' = T(s, a)$
- Receive outcome $R(s, a)$
- Get observation $o = O(s', a)$
- Update belief $b' = \tau(b, a, o)$
- Start over

$$V^\pi(b) = \sum_{t=0}^{\infty} \gamma^t r(b_t, a_t)$$

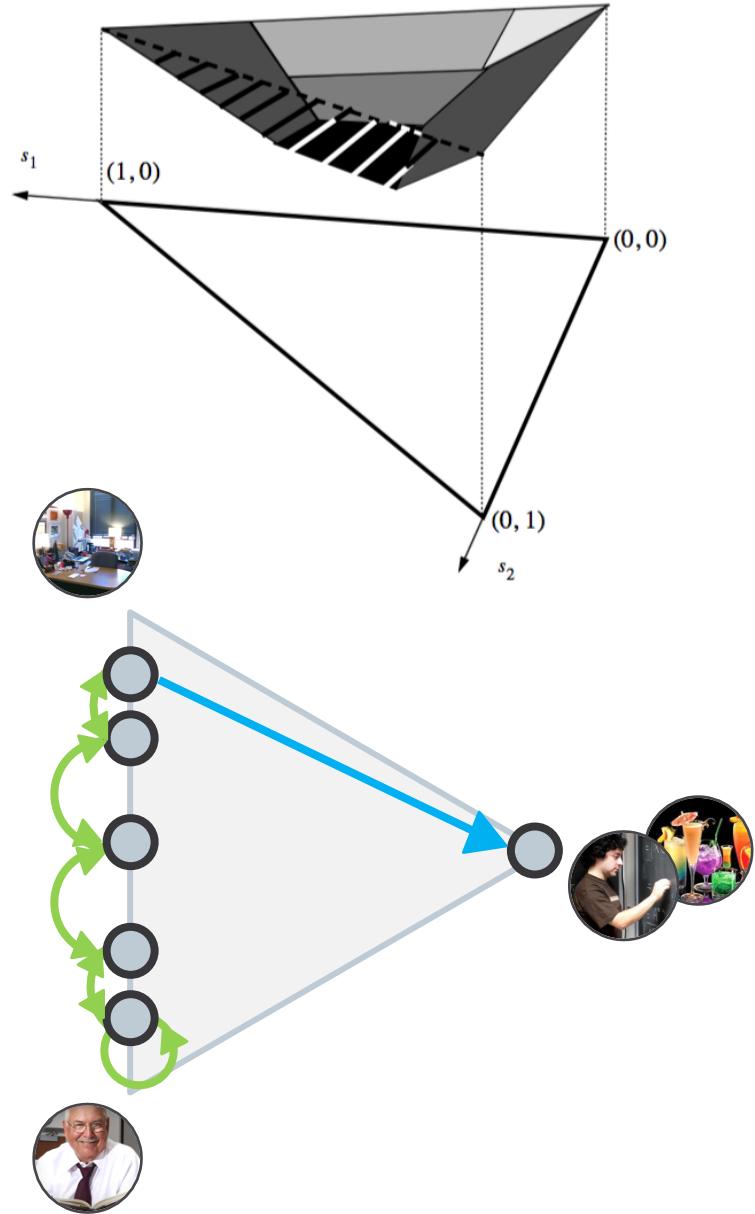
$$\pi^* = \operatorname{argmax}_{\pi} V^\pi$$



Resolution

The value function is always convex

- Certainty is preferable to uncertainty
- Gathering information is valuable



The solution can be discretized

- Optimal solution often visit a finite number of belief states
- The POMDP can then be reformulated as a (fully observable) MDP



Take home message

POMDPs allow to model:

- sequential decision making in a complex environment (MDP)
- subjectivity about the state of the world (PO)

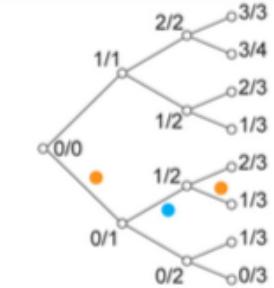
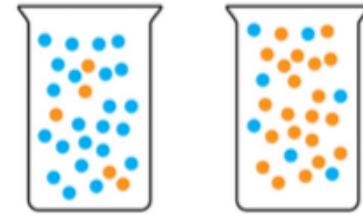
POMDPs can capture:

- information gathering as an economic decision
- irrational behaviour as an optimal policy based on wrong representations



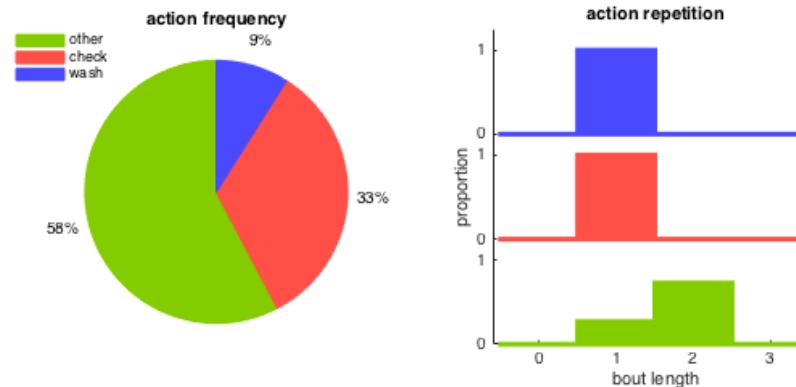
Perspectives

Information sequential sampling
with varying payoffs



[Averbeck 2015, PCB]

Errors as exploratory behaviour in
reversal learning tasks



Checking behaviour in OCD

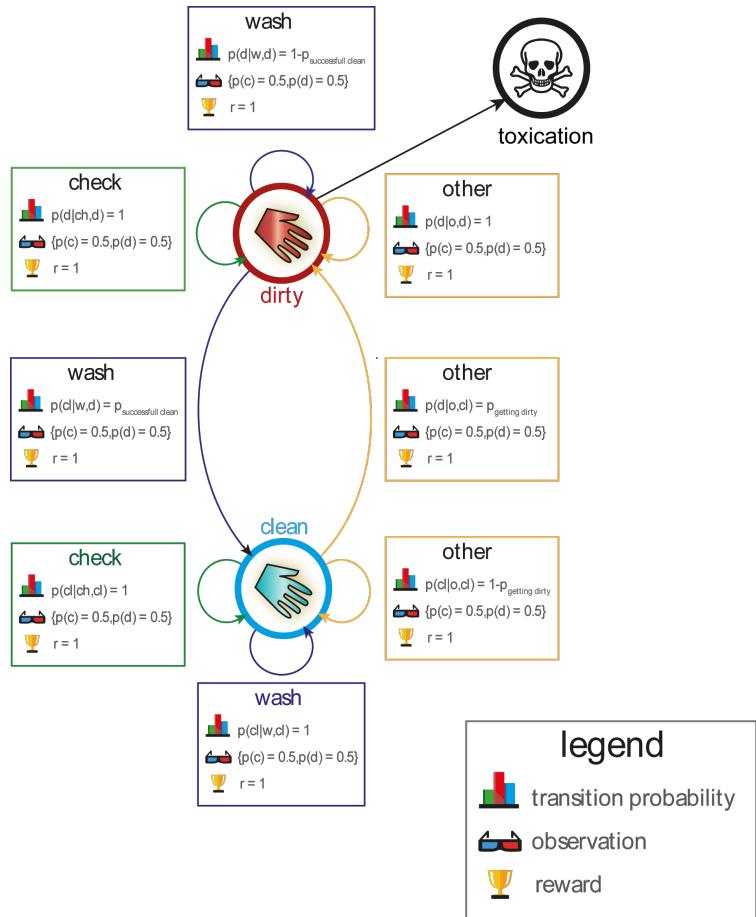


Questions?

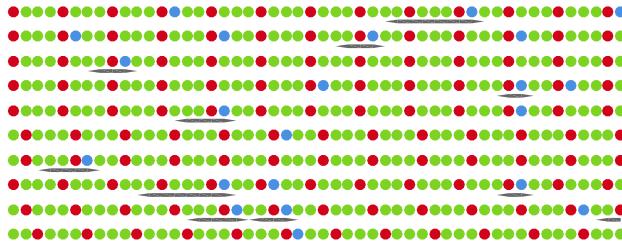
Thank you for your attention



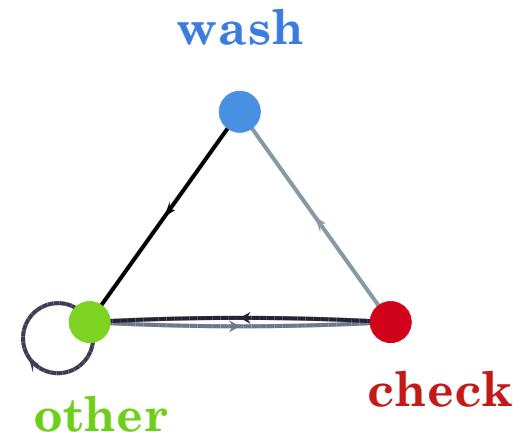
OCD model



optimal



wash



OCD

