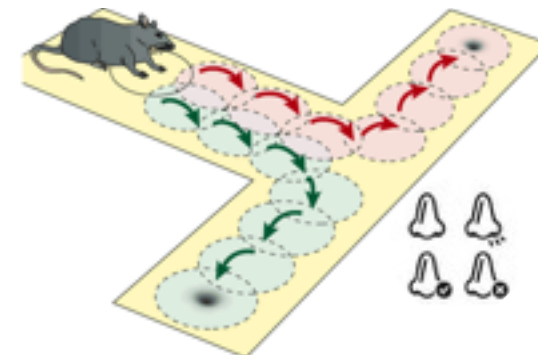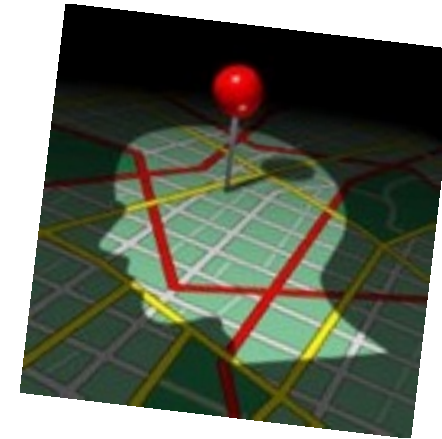# Partially Observable Markov Decision Processes

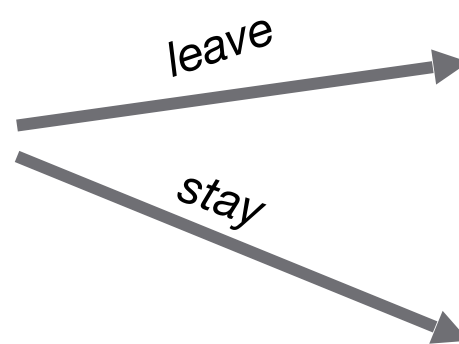Lionel Rigoux & Frederike Petzschner

# 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
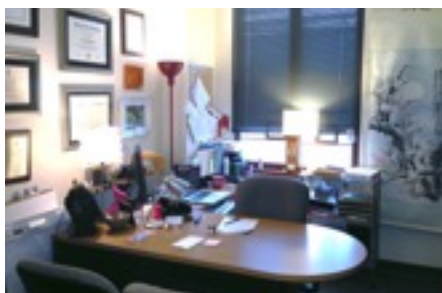    > reward location changes every trials, as cued by a smell
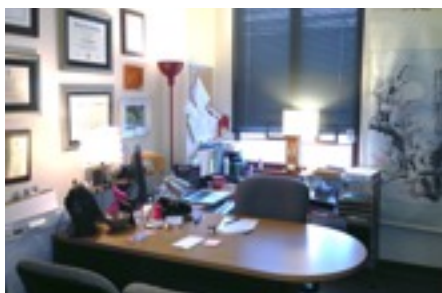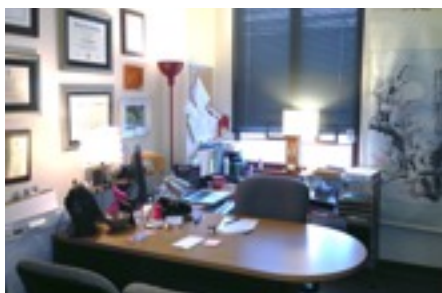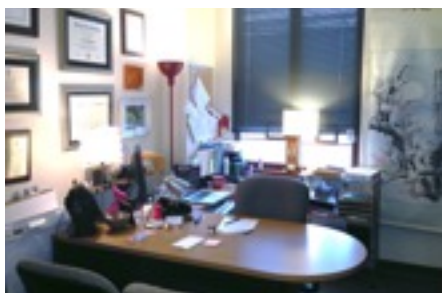
leave

leave

stay

*I'm gonna pretend to go to work now.*

leave

stay

I'm gonna pretend to go to work now.

leave

stay

I'm gonna pretend to go to work now.

leave

stay

stay

I'm gonna pretend to go to work now.

leave

stay

stay

leave

**state**



leave

stay

stay

leave

**state**　　　　**action**



leave

stay

stay

leave

# state          action          outcome

**state**          **action**                    **outcome**



leave                R = 100

stay                 R =   30

stay                 R =   30

leave                R = -40

**state**              **action**                          **outcome**



R = 100

R =   30

R = -40

**state**          **action**                              **outcome**



R = 100



?

R =   30



R = -40

**state**　　　　**action**　　　　　　　**outcome**



R = 100



*stay*

R =　30



R = -40

**state**　　　　**action**　　　　　　　　**outcome**



R = 100

*stay*

R =　30

R = -40

*leave*

**state**



$S_0$



$S_1$

**state**

*not known*



$S_0$



$S_1$

**state**

*not known*

**belief**

$b=p(s=S_1)$

$p(s=S_1) = 0$

$S_0$

$S_1$

$p(s=S_1) = 1$

**state**
*not known*

$S_0$

$S_1$

**belief**
$b=p(s=S_1)$

$p(s=S_1) = 0$

$p(s=S_1) = 1$

**actions and payoff function**

**state**

*not known*

**belief**

$b=p(s=S_1)$

**actions and payoff function**

$p(s=S_1) = 0$



$S_0$

$S_1$

$p(s=S_1) = 1$

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

**state**
*not known*

**belief**
$b=p(s=S_1)$

**actions and payoff function**

$p(s=S_1) = 0$

$S_0$

Expected reward

100

-40

0

belief

1

$S_1$

$p(s=S_1) = 1$

$E[R](a) = p(x=0) \, R_0(a) + p(x=1) \, R_1(a)$

**state**
*not known*

$S_0$

$S_1$

**belief**
$b=p(s=S_1)$

$p(s=S_1) = 0$

$p(s=S_1) = 1$

**actions and payoff function**

*stay*

Expected reward

100

-40

0    belief    1

$E[R](a) = p(x=0) \, R_0(a) + p(x=1) \, R_1(a)$

**state**
*not known*

$S_0$

$S_1$

**belief**
$b = p(s = S_1)$

$p(s = S_1) = 0$

$p(s = S_1) = 1$

**actions and payoff function**

*leave*

*stay*

Expected reward

100

-40

0    belief    1

$$E[R](a) = p(x=0) \, R_0(a) + p(x=1) \, R_1(a)$$

**state**

*not known*

$S_0$

$S_1$

**belief**

$b = p(s = S_1)$

$p(s = S_1) = 0$

$p(s = S_1) = 1$

**actions and payoff function**

*leave*

100

*stay*

Expected reward

-40

0

belief

1

$E[R](a) = p(x=0) \, R_0(a) + p(x=1) \, R_1(a)$

**state**

*not known*

$S_0$

$S_1$

**belief**

$b=p(s=S_1)$

$p(s=S_1) = 0$

$p(s=S_1) = 1$

**actions and payoff function**

*leave*

*stay*

100

-40

Expected reward

0

1

belief

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

**state**
*not known*

$S_0$

$S_1$

**belief**
$b = p(s = S_1)$

$p(s = S_1) = 0$

$p(s = S_1) = 1$

**actions and payoff function**

*leave*

*stay*

Expected reward

100

-40

0     belief     1

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

**state**
*not known*

$S_0$

$S_1$

**belief**
$b=p(s=S_1)$

$p(s=S_1) = 0$

$p(s=S_1) = 1$

**actions and payoff function**

*leave*

*stay*

Expected reward

100

-40

0

belief

1

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

# state
*not known*



$S_0$



$S_1$

# belief
*b=p(s=S$_1$)*

$p(s=S_1) = 0$

$p(s=S_1) = 1$

# actions and payoff function



*leave*

*stay*

Expected reward

100

-40

0

belief

1

$$E[R](a) = p(x=0) \, R_0(a) + p(x=1) \, R_1(a)$$

# state

*not known*



$S_0$



$S_1$

# belief

$b=p(s=S_1)$

$p(s=S_1) = 0$

$p(s=S_1) = 1$

# actions and payoff function

**optimal policy**

*leave*

$R(leave, s=1) = 5$

*stay*

Expected reward

100

-40

0      **belief**      1

$$E[R](a) = p(x=0)\, R_0(a) + p(x=1)\, R_1(a)$$

**state**

$S_0$

$S_1$

**observation function**

*provide information about state*

**belief**

$p(s=S_1) = 0$

$p(s=S_1) = 1$

**state**

$S_0$

$S_1$

**observation function**

*provide information about state*

| | leave | stay | listen |
|---|---|---|---|
| **noises** | 0 | 0.5 | 0.15 |
| **no one** | 1 | 0.5 | 0.85 |

| | leave | stay | listen |
|---|---|---|---|
| **noises** | 1 | 0.5 | 0.85 |
| **no one** | 0 | 0.5 | 0.15 |

**belief**

$p(s=S_1) = 0$

$p(s=S_1) = 1$

**state**

**observation function**

*provide information about state*

**belief**

$p(s=S_1) = 0$


$S_0$

| | leave | stay | listen |
|---|---|---|---|
| noises | 0 | 0.5 | 0.15 |
| no one | 1 | 0.5 | 0.85 |

$$b' \sim p(o|s',a) \sum_{s} p(s'|s,a)b(s)$$


$S_1$

| | leave | stay | listen |
|---|---|---|---|
| noises | 1 | 0.5 | 0.85 |
| no one | 0 | 0.5 | 0.15 |

$p(s=S_1) = 1$

# state

**state**

$S_0$

$S_1$

# observation function

*provide information about state*

| | leave | stay | listen |
|---|---|---|---|
| noises | 0 | 0.5 | 0.15 |
| no one | 1 | 0.5 | 0.85 |

$$b' \sim p(o|s',a) \sum_s p(s'|s,a)b(s)$$

| | leave | stay | listen |
|---|---|---|---|
| noises | 1 | 0.5 | 0.85 |
| no one | 0 | 0.5 | 0.15 |

# belief

$p(s=S_1) = 0$

$p(s=S_1) = 1$

**state**

$S_0$

$S_1$

**observation function**

*provide information about state*

| | leave | stay | listen |
|---|---|---|---|
| noises | 0 | 0.5 | 0.15 |
| no one | 1 | 0.5 | 0.85 |

$$b' \sim p(o|s',a) \sum_s p(s'|s,a)b(s)$$

| | leave | stay | listen |
|---|---|---|---|
| noises | 1 | 0.5 | 0.85 |
| no one | 0 | 0.5 | 0.15 |

**belief**

$p(s=S_1) = 0$

*listen*
no one

$p(s=S_1) = 1$

**state**

**observation function**
*provide information about state*

**belief**

$p(s=S_1) = 0$

| | leave | stay | listen |
|---|---|---|---|
| noises | 0 | 0.5 | 0.15 |
| no one | 1 | 0.5 | 0.85 |

$S_0$

$$b' \sim p(o|s',a) \sum_s p(s'|s,a)b(s)$$

*listen*
no one

| | leave | stay | listen |
|---|---|---|---|
| noises | 1 | 0.5 | 0.85 |
| no one | 0 | 0.5 | 0.15 |

$S_1$

$p(s=S_1) = 1$

**state**

**observation function**

*provide information about state*

**belief**

$p(s=S_1) = 0$

$S_0$

| | *leave* | *stay* | *listen* |
|---|---|---|---|
| **noises** | 0 | 0.5 | 0.15 |
| **no one** | 1 | 0.5 | 0.85 |

$$b' \sim p(o|s',a) \sum_s p(s'|s,a)b(s)$$

| | *leave* | *stay* | *listen* |
|---|---|---|---|
| **noises** | 1 | 0.5 | 0.85 |
| **no one** | 0 | 0.5 | 0.15 |

$S_1$

*listen*
no one

*listen*
no one

$p(s=S_1) = 1$

# state

**observation function**

*provide information about state*

# belief

$p(s=S_1) = 0$

**S_0**

| | leave | stay | listen |
|---|---|---|---|
| **noises** | 0 | 0.5 | 0.15 |
| **no one** | 1 | 0.5 | 0.85 |

$$b' \sim p(o|s',a) \sum_s p(s'|s,a)b(s)$$

| | leave | stay | listen |
|---|---|---|---|
| **noises** | 1 | 0.5 | 0.85 |
| **no one** | 0 | 0.5 | 0.15 |

**S_1**

*listen*
no one

*listen*
no one

$p(s=S_1) = 1$

## state



$S_0$



$S_1$

## observation function

*provide information about state*

|  | leave | stay | listen |
|---|---|---|---|
| **noises** | 0 | 0.5 | 0.15 |
| **no one** | 1 | 0.5 | 0.85 |

$$b' \sim p(o|s',a) \sum_s p(s'|s,a)b(s)$$

|  | leave | stay | listen |
|---|---|---|---|
| **noises** | 1 | 0.5 | 0.85 |
| **no one** | 0 | 0.5 | 0.15 |

## belief

$p(s=S_1) = 0$

*leave*
noises

*listen*
no one

*listen*
no one

$p(s=S_1) = 1$

# state

**observation function**

*provide information about state*

# belief

$p(s=S_1) = 0$


S0

| | leave | stay | listen |
|---|---|---|---|
| noises | 0 | 0.5 | 0.15 |
| no one | 1 | 0.5 | 0.85 |

$$b' \sim p(o|s', a) \sum_s p(s'|s, a)b(s)$$

| | leave | stay | listen |
|---|---|---|---|
| noises | 1 | 0.5 | 0.85 |
| no one | 0 | 0.5 | 0.15 |


S1

*leave*
noises

*listen*
no one

*listen*
no one

$p(s=S_1) = 1$

## state space



*stay*    *listen*

0.9     1

*leave*

1

$S_0$

0.1

$S_2$

1

0.9     1

$S_1$

## belief space



$s_1$    $(1,0)$

$(0,0)$

$(0,1)$

$s_2$

# state space

# belief space



*stay*

*listen*

0.9

1

*leave*

1

0.1

**S₀**

0.9

**S₁**

1

1

**S₂**

1

*Outcomes and observations are determined by the real state space*

*Policy relies on the belief state*

$s_1$  (1,0)

(0,0)

(0,1)

$s_2$

# POMDP Formalism

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

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

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

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

# POMDP Formalism

**MDP**
- $S$    set of states
- $A$    set of actions
- $T$    transition matrix   $S{\times}A \to S$
- $R$    reward function   $S{\times}A \to \mathbb{R}$
- $\gamma$    discount factor

**POMDP extension**
- $\Omega$    set of observations
- $O$    observation probabilities $S{\times}A{\times}\Omega \to [0,1]$
- $B$    belief space
- $r$    reward function   $B{\times}A \to \mathbb{R}$
- $\tau$    belief update function   $B{\times}A{\times}\Omega \to B$

**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

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

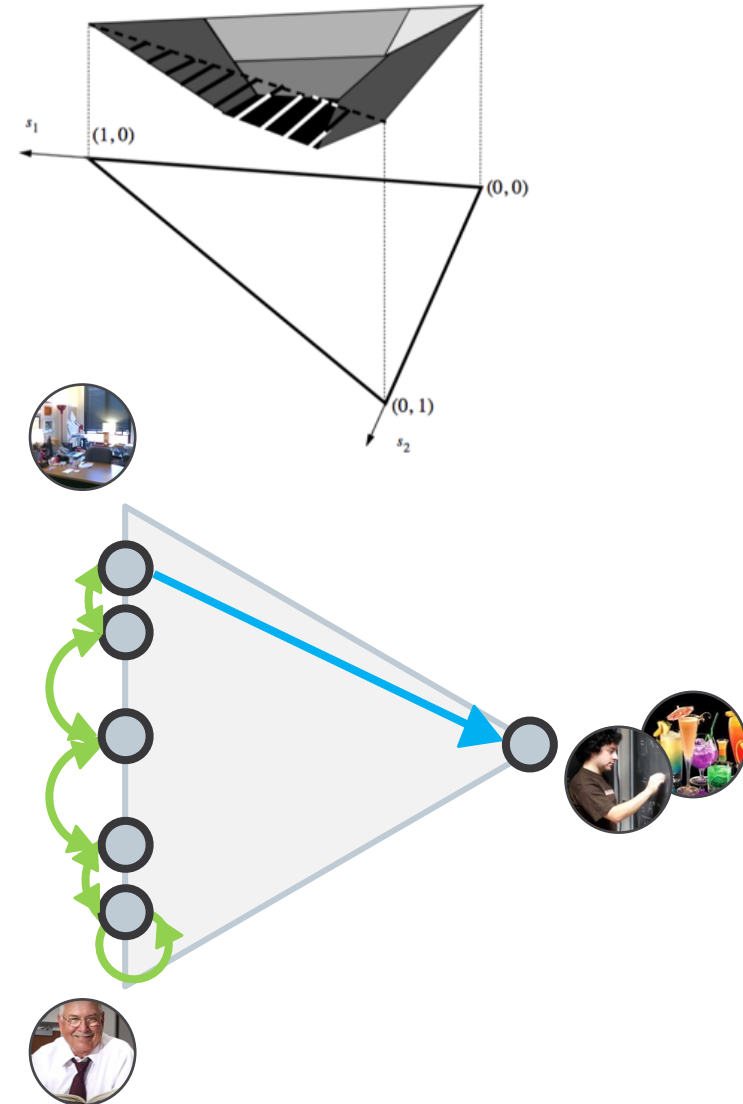$$\pi^{\star} = \operatorname*{argmax}_{\pi} V^{\pi}$$

# 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, evolving environment (MDP)
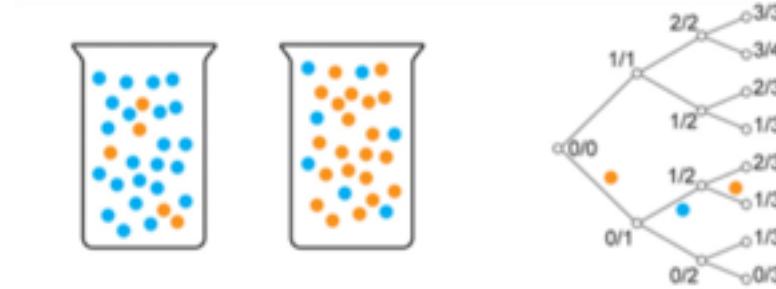
- subjectivity about the state of the world (PO)

POMDPs can capture:

- information gathering as an economic decision

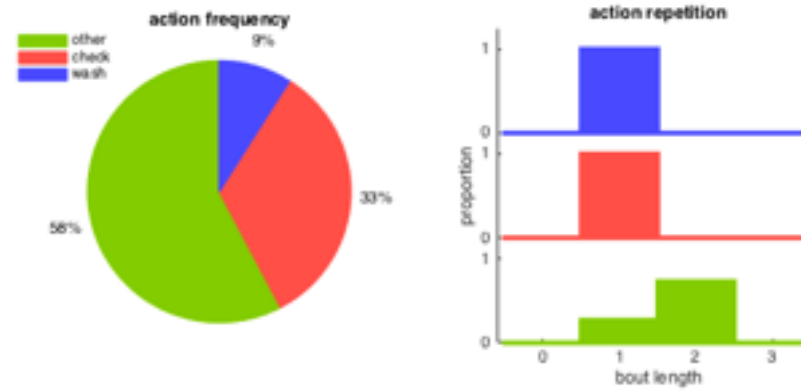- irrational behaviour as an optimal policy based on wrong representations

# Perspectives

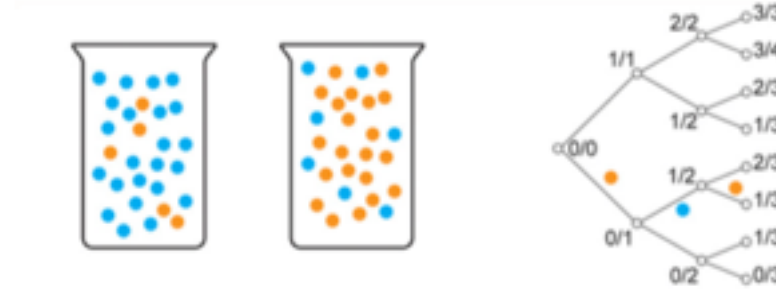Information sequential sampling with varying payoffs

Errors as exploratory behaviour in reversal learning tasks

Checking behaviours in OCD

# Perspectives

Information sequential sampling
with varying payoffs

Errors as exploratory behaviour in
reversal learning tasks

Checking behaviours in OCD

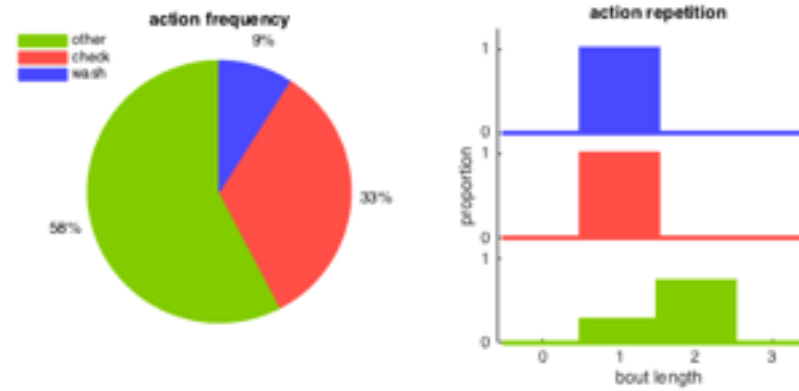[Averbeck 2015, PCB]

# Questions?

Thank you for your attention