

Last Time:

- Intro to HRL

This Time:

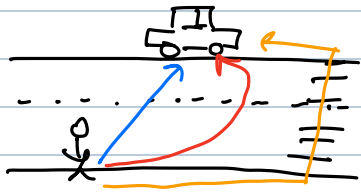
- sequential decision-making
- MDPs

Lecture 2

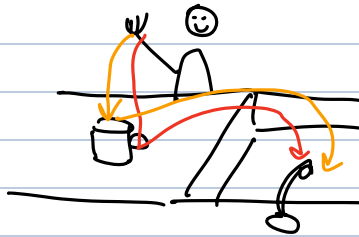
HRL, Fall '24

Andrea Bajcsy

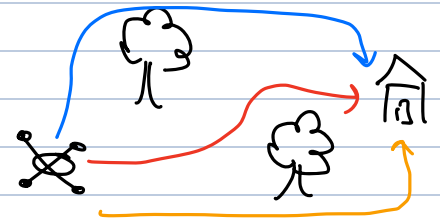
Sequential decision-making is everywhere - play games, making life decisions-- and its present in interaction. In HRI, sequential decision-making will form the "mathematical backbone" of how we model people, robots, and their interaction.



how will this person move to their car?



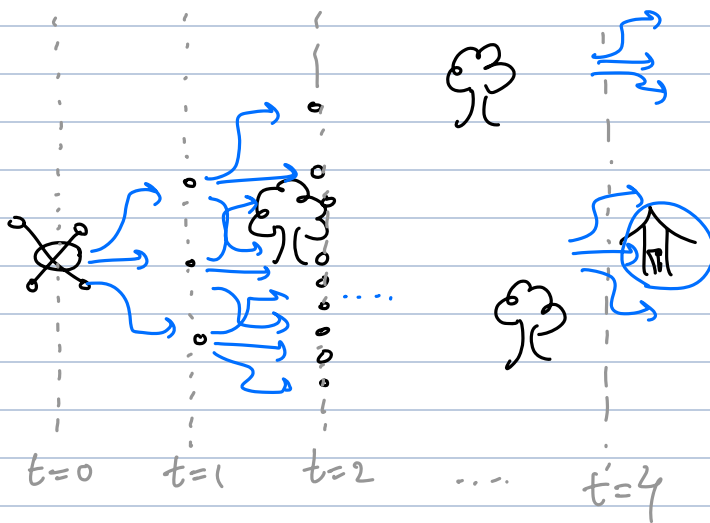
how will this person clean their mug?



how will drone fly home?

ⓐ What makes sequential decision-making hard?

ⓐ naive solution grows exponentially in time horizon.



action space
 $|A| = 3$

start from ~~X~~ state

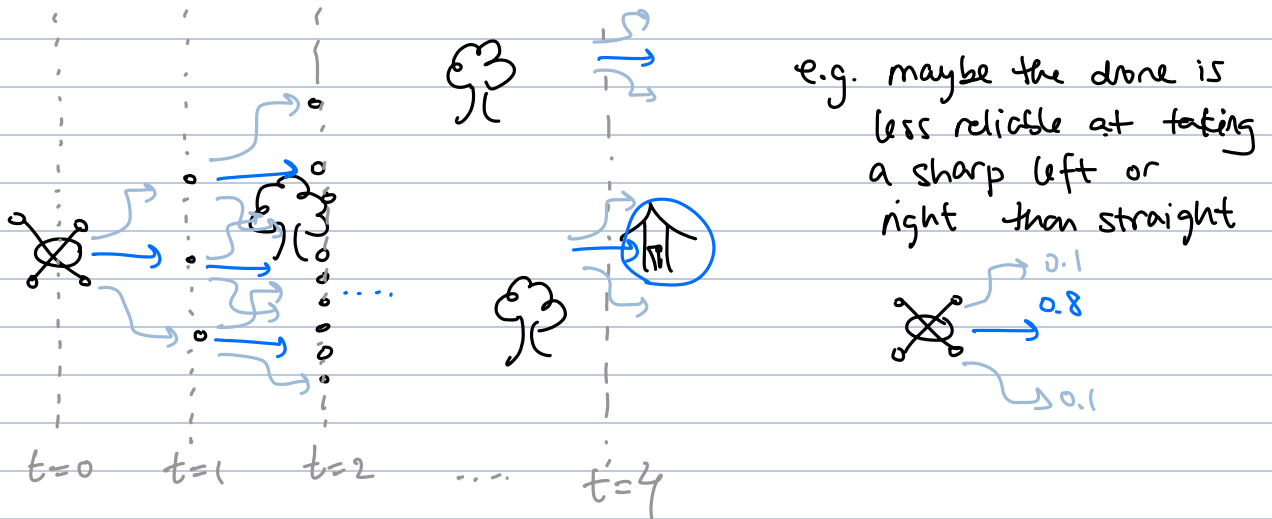
there are:

$$|A|^{T+1}$$

sequences of decisions to choose from.

$= 3^4 = 81$ possible sequences of decisions w/ just 4 seconds ∞

1A2 Outcomes of taking actions can sometimes be stochastic.

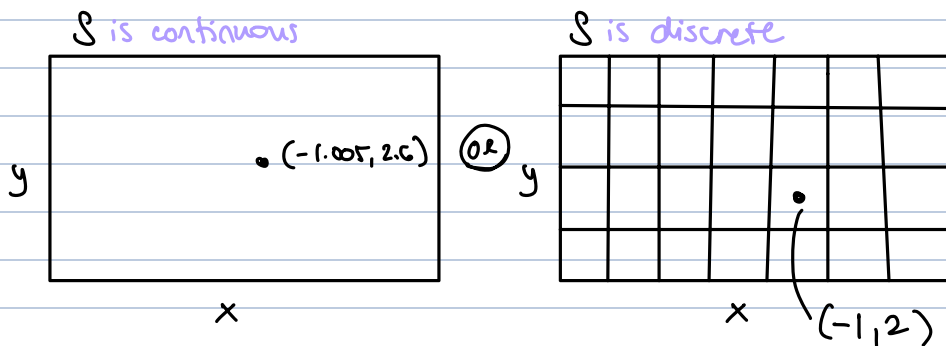


Markov Decision Processes:

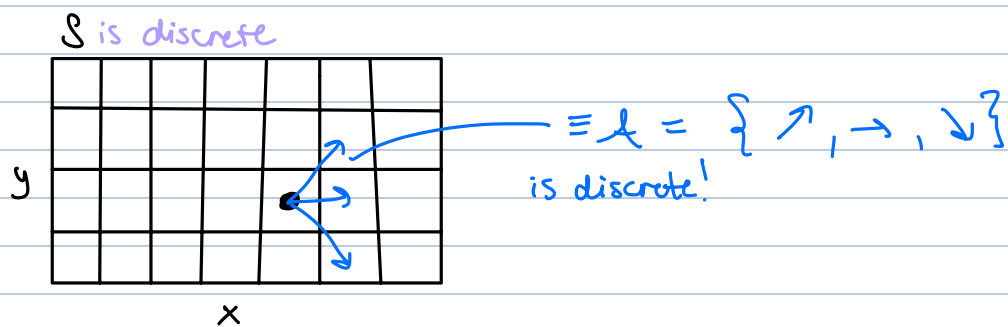
this is a mathematical model for sequential decision-making in a fully observable, stochastic environment with a **Markovian transition** model. Let's break down this model into the key "ingredients" and modeling assumptions:

"MDP is a tuple of $\langle S, A, T, r \rangle$ "

- $s \in \underline{S}$ = state space.

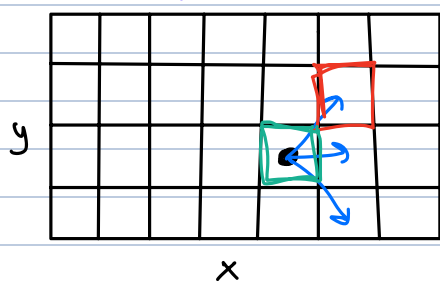


- $a \in \underline{A}$ = action space. This is what our agent can do.



$T: S \times A \rightarrow S$ is the transition function. Is just a map
 $T: S \times A \rightarrow \Delta(S)$ distribution over next states
 ↪ stochastic transition function

from the space of states S and actions to a (distribution) over next states.



$$T(s_{t+1} | s_t, a_t)$$

$$P(s_{t+1} | s_t, a_t)$$

$$P(s_{t+1} = \square | s_t = \square, a_t = \nearrow) = 0.9$$

⊗ In the transition function is precisely where the "Markov" part of Markov Decision-Making Processes comes into play.

Specifically, the Markov Property states that the future is independent of the past, given the present:

$$P(\text{future} | \text{present, past}) = P(\text{future} | \text{present})$$

In MDPs, this translates to the outcome of actions only depending on the current state, and not a history:

$$P(s_{t+1} | s_t, a_t, s_{t-1}, a_{t-1}, \dots, s_0) \stackrel{\uparrow}{=} P(s_{t+1} | s_t, a_t)$$

by Markov Property!

• $r: S \rightarrow \mathbb{R}$ is the reward function. Its called "instantaneous reward" b/c its only thinking about where you are RIGHT NOW!

reward ↴

○	○	○	○	○	○	○
○	○	○	↻ -10	○	○	↕ +10
○	↻ -10	○	●	○	○	○
○	○	○	○	↻ -10	○	○

x

we have all these components of what it means to make decisions, but we need a way for our agent to solve or know what the "best" decision is for any state they may be at.

The object we seek to solve for is a policy.

$\pi: S \rightarrow \mathcal{A}$ is a mapping from states to actions

$\pi: S \rightarrow \mathcal{A}$

→	→	→	→	→	→	→
→	→	→	↻	→	→	↕
→	↻	→	→	→	→	→
→	→	→	↻	→	→	→

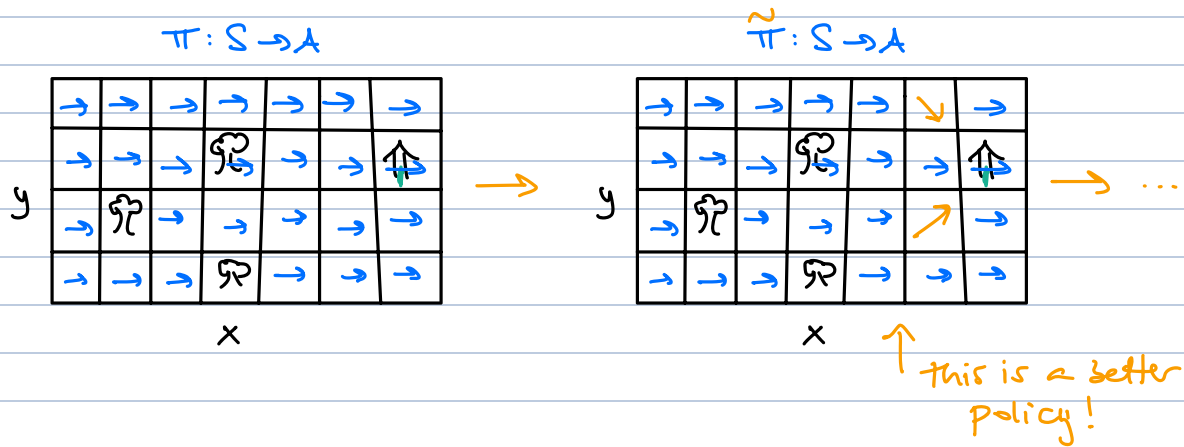
x

Evaluate the quality of a policy π by the expected cumulative reward induced by the policy.

$$R(s_0, s_1, s_2, \dots) = r(s_0) + r(s_1) + r(s_2) + \dots$$

$$= \sum_{t=0}^{\infty} r(s_t)$$

An optimal policy $\pi^*: S \rightarrow A$ yields the highest expected reward for all states.



• $\gamma =$ discount factor, $\gamma \in [0, 1]$

$$R(s_0, s_1, s_2, \dots) = \gamma^0 r(s_0) + \gamma^1 r(s_1) + \gamma^2 r(s_2) + \dots$$

$$= \sum_{t=0}^{\infty} \gamma^t r(s_t)$$

Discount factor describes the preference of an agent for current rewards over future ones when $\gamma < 1$. When $\gamma = 1$, then our agent wants the max reward over all time steps.

⊕ discounting appears to be a good model of human + animal preferences over time.

Ultimately, we are searching for a policy $\pi^*: S \rightarrow A$ that maximizes the sum of discounted rewards

$$\pi^* := \arg \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right]$$

$$s_{t+1} \sim P(s_{t+1} | s_t, \pi)$$

MORE INFO:

"Artificial Intelligence: A Modern Approach" by Russel & Norvig.