

Last Time:

□ Intro to HRL

lecture 2

HRL, FALL '25

Andrea Bajcsy

This Time:

□ sequential decision-making

□ MDPs

□ Policies, Values, Bellman Eqn.

What is sequential decision-making?

Sequential decision-making involves making a sequence of decisions over time, where each decision affects future outcomes and decisions.

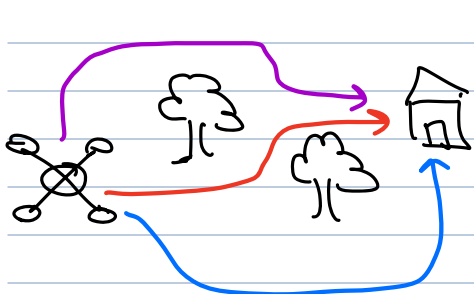
↳ different from one-shot decisions which are a self-contained single decision (e.g. img. classification)

Why do we care?

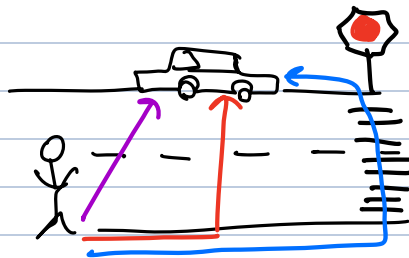
Sequential decision-making is everywhere - playing games, planning a trip or career - and its present in interaction!

In HRI, sequential decision-making will form the "mathematical backbone" of how we model people, robots, and their interaction:

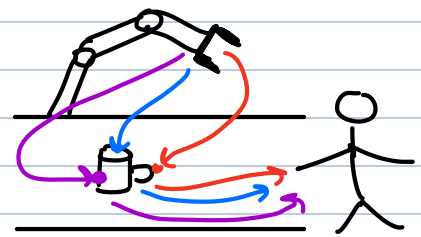
- robots must plan ahead, considering long-term consequences
- humans are sequential decision-makers too: must model them well
- interaction is an ongoing exchange where agents influence each other



How will R fly home?



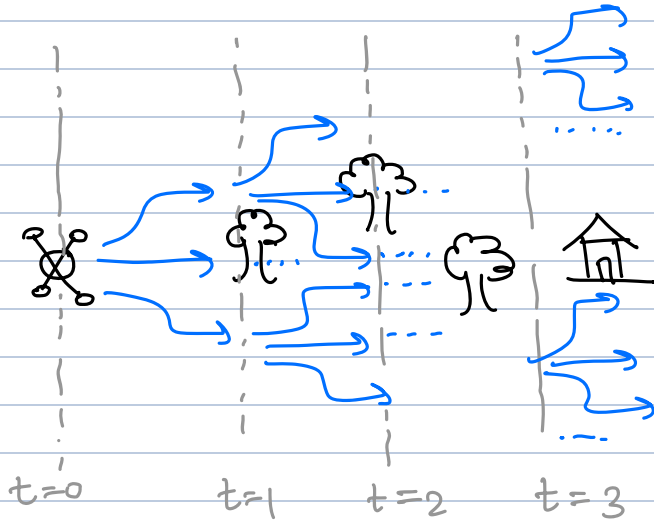
How will H go to their car?



How should R hand H the mug?

Q Why is it hard?

A1 decision tree grows exponentially in time horizon.
↳ too many possibilities to compute



$$|A| = 3 \Rightarrow$$

starting from ~~X~~ state
there are

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

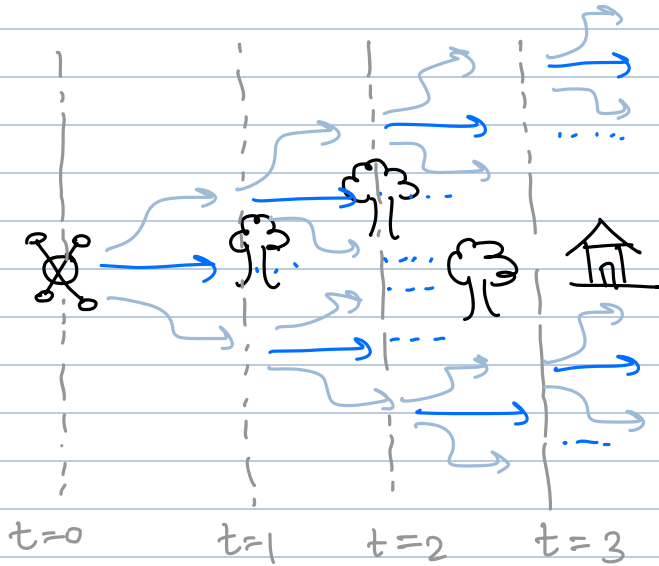
sequences of decisions

$$= 3^4 = 81 \text{ sequences}$$

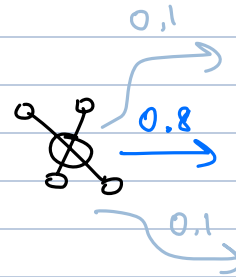
with time-horizon of 3! Δ

⇒ Naively planning quickly becomes impossible!

A2 outcomes of taking actions can be stochastic
↳ same action doesn't always have the same result.



e.g. R chooses to go straight,
but wind might push it
left or right.



⇒ Even with same action sequence,
outcomes may not be the same!

A3 Hard to assign credit to the right past action(s)
↳ delayed rewards make it unclear which action was responsible.



✓ successfully home, but was it all good...?
X fail to get home, but was it all bad...?

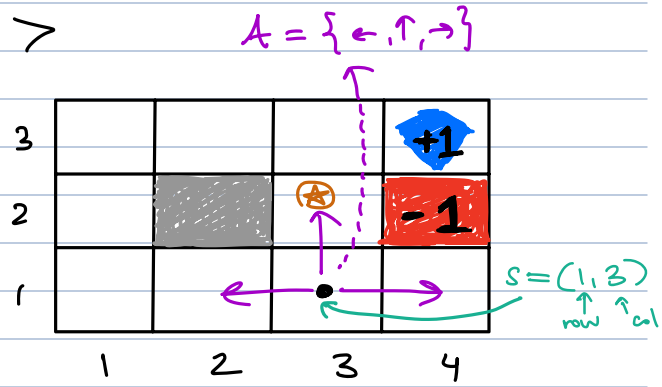
⇒ If rewards come late, hard to tell which
actions were good/bad!

Markov Decision Processes (MDPs)

MDPs are a mathematical model for sequential decision-making in a fully observable, stochastic, discrete-time environment with a markovian transition model.

$$M = \langle S, A, T, r \rangle$$

- states: $s \in S$ ↗ state space
- actions: $a \in A$ ↗ action space
- transition function: $T: S \times A \rightarrow S$



$$T(s, a, s') \quad \underline{\underline{\text{OR}}}$$

$$P(s' | s, a) \quad \underline{\text{or}}$$

$s' \sim P(\cdot | s, a)$ "is distributed as"

$$T(s=(1,3), a=\uparrow, s'=(2,3)) = 0.8$$

$$P(s' = (2, 3) \mid s = (1, 3), a = \uparrow) = 0.8$$

$$s' \sim P(\cdot | s=(1,3), a=\uparrow) = \begin{cases} (2,3) \rightarrow 0.8 \\ (1,2) \rightarrow 0.1 \\ (1,4) \rightarrow 0.1 \end{cases}$$

- reward function : $r : S \times A \times S \rightarrow \mathbb{R}$

$$\frac{r(s, a, s')}{r(s, a)} \stackrel{\text{or}}{=} \frac{r(s')}{1}$$

⑤ $r(s=(2,3), a=\rightarrow, s'=(2,4)) = -1$

$$r(s' = (2, 4)) = -1$$

$$r(s' = (3, 4)) = +1$$

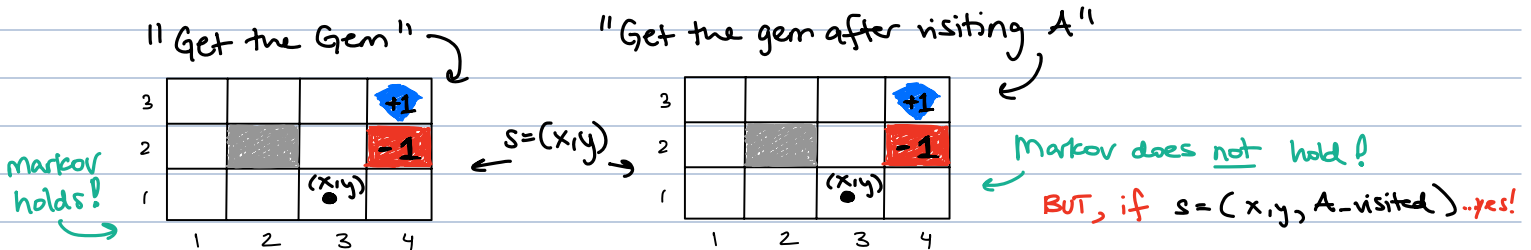
② what is "markov" about this?

The transition function is where the markov part comes into play.

- Markov Property: the future is independent of the past, given the present.

This means that action outcomes depend only on the current state, not history!

$$P(s_{t+1} | s_t, a_t, s_{t-1}, a_{t-1}, \dots, s_0) = P(s_{t+1} | s_t, a_t) \quad \text{! makes planning easier}$$



! If a problem is Markov depends on how we define the state

Policies

The object we ultimately want to solve for is the policy: OR, knowing what is the best decision to make @ any state.

$\pi: S \rightarrow A$
"deterministic, $\pi(s)$ "

$\pi: S \rightarrow \Delta A$
"stochastic: $\pi(a|s)$ "
probability simplex, "distribution over actions"

ex. policy (not very good one tho...)

$\pi: S \rightarrow A$

3	→	→	→	+1
2	→		→	-1
1	→	→	→	→
	1	2	3	4

Evaluate the quality of a policy π by the expected cumulative reward induced by the policy. \otimes we will get to the "expected" bit in a bit...

① "roll out" or simulate π from each $s_0 \in S$ to get a trajectory sequence $\{s_0, \pi(s_0), s_1, \pi(s_1), \dots\}$
 a_0 a_1

② Evaluate the utility of the rollout, called cumulative reward:

$$R(s_0, s_1, s_2, \dots) \triangleq r(s_0) + r(s_1) + \dots = \sum_{t=0}^{\infty} r(s_t) \quad \left[\begin{array}{l} \text{or } r(s_t, a_t) \\ r(s_t, a_t, s_{t+1}) \end{array} \right]$$

$\pi_1: S \rightarrow A$

3	→	→	→	+1
2	→		→	-1
1	→	→	→	→
	1	2	3	4

$\pi_2: S \rightarrow A$

3	→	→	→	+1
2	↑		↑	-1
1	↑	→	↑	←
	1	2	3	4

← better policy b/c more reward!

• an optimal policy ($\pi^*: S \rightarrow A$) yields the highest cumulative reward

⚠ But what counts as a "good" policy depends on the reward function!

Sparse vs. Dense Rewards

Problem: If the only reward is a +1 at the goal (and -1 @ lava), many policies look equally good because there is no penalty for inefficiency! In other words, sparse rewards often underspecify what we "actually" want the agent to do!

⊗ Underspecified rewards is a HUGE problem in HRL and AI Alignment!!!

Solution: Dense rewards — frequent feedback (good / bad) to the agent over time, instead of just rewarding outcomes at the end.

- ⊕ helps guide the agent to behave in the way we actually want.
- ⊕ specify not just what the agent should do (i.e. goal), but how to do it.
- ⊕ faster agent learning
- ⊖ HARD TO SPECIFY + PROVIDE → see later lectures on this!

Example: "Living" ^{→ timestep} rewards are a kind of dense reward that incentivizes efficiency by penalizing necessary steps.

π^*

3	→	→	→	+1
2	↑		←	-1
1	↑	←	←	↓
	1	2	3	4

$r(s) = -0.01$

↑ small living penalty

⇒ detours around obstacle + lane, i.e. "bumps" away from walls / obstacle

π^*

3	→	→	→	+1
2	↑		↑	-1
1	↑	→	↑	←
	1	2	3	4

$r(s) = -0.4$

↑ med. living penalty

⇒ no bumping against walls, direct path

π^*

3	→	→	→	+1
2	↑		→	-1
1	→	→	→	↑
	1	2	3	4

$r(s) = -2.0$

↑ large living penalty

⇒ This one wants to exit ASAP!

Discount Factor

- denoted by $\gamma \in [0, 1]$ ⇒ sometimes MDP tuple is written as $\langle S, A, T, r, \gamma \rangle$
- describes an agent's preference for current rewards over future ones when $\gamma < 1$. If $\gamma \equiv 1$, our agent wants max reward over all tsteps.

$$R(s_0, s_1, s_2, \dots) = \gamma^0 r(s_0) + \gamma^1 r(s_1) + \gamma^2 r(s_2) + \dots = \sum_{t=0}^{\infty} \gamma^t r(s_t)$$

↙ discounted cumulative reward ↘

⊗ discounting is a good model of human & animal preferences

↳ See: "models of Temporal Discounting 1937-2000" by Till Grüne-Yanoff.

→ connections w/ psych + econ.

Ultimately, we are searching for a policy $\pi^*: S \rightarrow A$ that maximizes the expected discounted cumulative reward:

$$\pi^* = \operatorname{argmax}_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$

"return the policy which maximizes..." π
 trade off shorter vs. longer horizon rewards γ^t
 rollouts are stochastic so take expectation over all trajectories under the policy π !
 let $\tau := (s_0, a_0, s_1, a_1, \dots)$
 $p_{\pi}(\tau) = p(s_0) \prod_{t=0}^{\infty} p(s_{t+1} | s_t, a_t)$
 * note: this is for deterministic π ! for stochastic we have
 $= p(s_0) \prod_t \pi(a_t | s_t) p(s_{t+1} | s_t, a_t)$