

Last Time:

□ Bellman Eqn.

□ Value Iteration + RL

Lecture 4

421, FALL '25

Andrea Bajcsy

This Time:

□ RL summary

□ POMDPs

→ HW #1 has been released on Canvas! Due: Sept. 18

→ Yilin's OHs: Mondays, 11am - 12pm in NSH 4306

Reinforcement learning (RL)

[Q] In MDPs, we assumed we know $P(s'|s, a)$ and rewards $r(s, a)$... but what if we don't?

In reality, it's hard to know P and r explicitly, so how do we obtain π^* ?

[A] RL! It allows an agent to learn optimal policies by interacting with the environment.

⇒ agent tries out actions, observes rewards, and updates their policy to maximize rewards.

Model-based vs. Model-free RL

"model of the world"

↳ High-level diff if the agent has access to (or learns) a transition function + rewards

(A) Model-based: 1) collect data from environment via some π

$$\mathcal{D} := \{(s_t, a_t, s_{t+1}, r_t)_{t=1}^N\}$$

2) use supervised learning to fit empirical MDP model:

↳ count s' for each s, a ; normalise; get $\hat{P}(s'|s, a)$

↳ discover each $\hat{r}(s, a, s')$ when we experience (s, a, s')

3) solve MDP via techniques above (+more!)

(B) Model-free: 1) collect data from environment via some π

$$\mathcal{D} := \{(s_t, a_t, s_{t+1}, r_t)_{t=1}^N\}$$

2) Directly learn a Q-function or π via trial + error
"Q-learning" "direct policy search"

⇒ for more depth on topic see:

- Sutton & Barto, "Reinforcement Learning."
- Russell & Norvig, "AI: A Modern Approach". Ch. 21

Partially-observable Markov Decision Processes (POMDPs)

POMDPs extend MDPs to settings where the agent cannot directly observe the true state. Instead, the agent must infer the state from noisy observations.

Formally, a POMDP is a tuple

$$\mathcal{M} = \langle S, A, T, r, \underbrace{\mathcal{O}}_{\text{new!}}, Z, \gamma \rangle$$

- states: $s \in S$ (partially observable)

- actions: $a \in A$

- transition : $T: S \times A \rightarrow S$
function $T(s, a, s') \propto P(s' | s, a)$

- reward : $r: S \times A \times S \rightarrow \mathbb{R}$
function $r(s, a, s') \propto r(s, a), r(s)$

- observations: $o \in \mathcal{O} \rightarrow$ observation space

- observation: $Z: S \times A \rightarrow \mathcal{O}$
function $Z(s', a, o') \propto$

$$P(o' | s', a) \propto$$

$$Z(s, o) = P(o | s)$$

observation (1) resulting state (since the world depends on determines what can be observed)

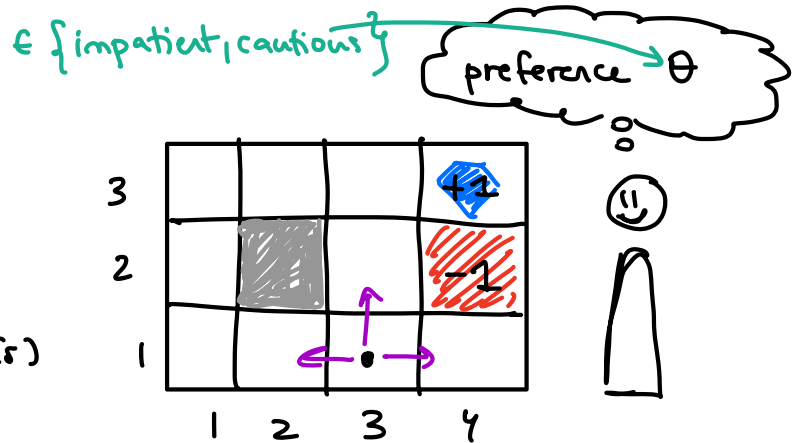
(2) previous action (since actions affect what can be perceived)

- discount $\gamma \in [0, 1]$

Why POMDPs? Many real-world problems have partial observability

→ ex. self-driving car doesn't know intentions of pedestrians & must infer them from observations $o \in \mathcal{O}$.

A key concept in POMDPs is the belief state. Since the R doesn't know s , it can maintain a belief state $b(s)$ as an estimate of state! It updates this belief over time via Bayes' Rule after taking actions + getting observations.

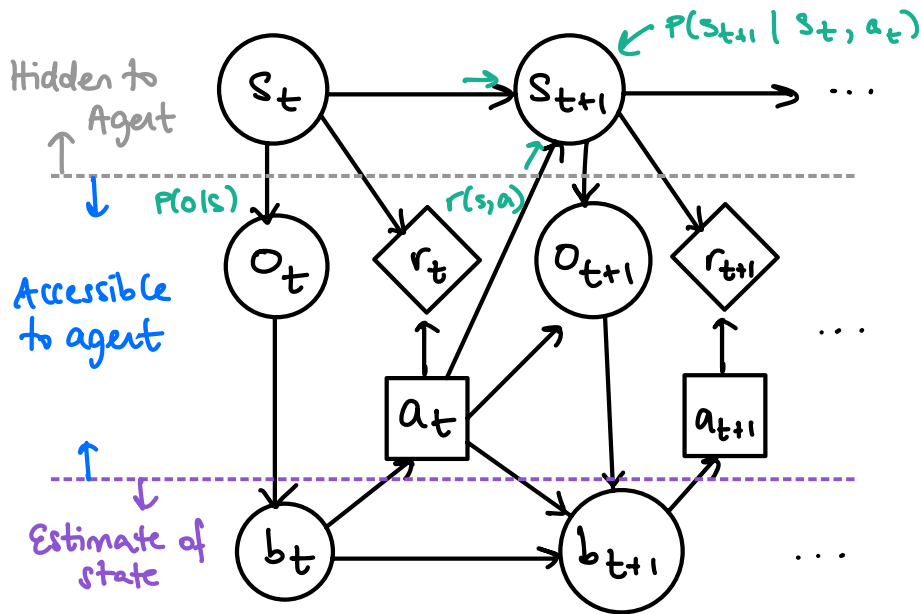


$$s := (x, y, \theta)$$

$$b(s) = P(s)$$

Belief state: $b(s) = P(s)$

a probability distribution over world states



ASIDE: graphical notation from decision networks / influence diagrams

○ = chance node (i.e. random var).

◇ = "utility" node

□ = decision node