Last Time

   ☐ Intro to Embodied AI Safety

This Time

   ☐ sequential decision-making review
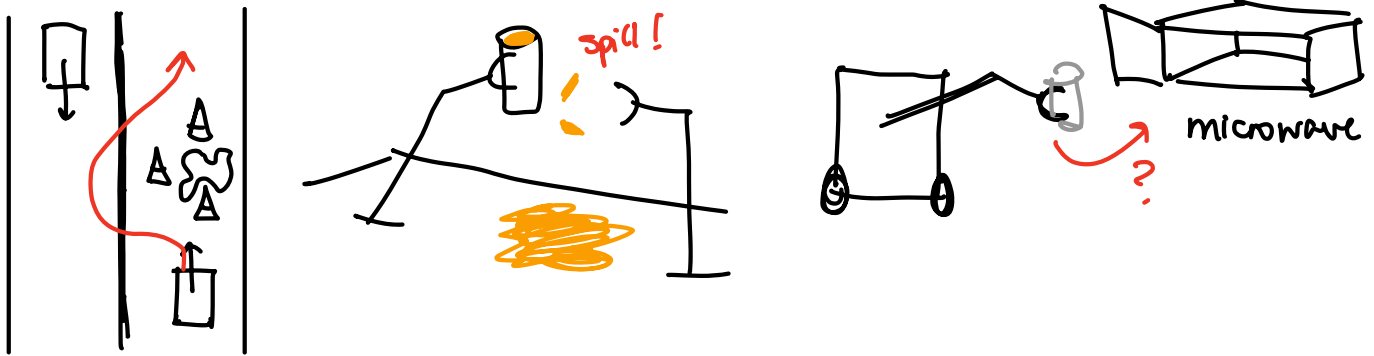
EAIS SP'26

Andrea Bajcsy

Sequential Decision-Making — b/c the first half of the class is all about safety w.r.t decisions.

e.g. I want an embodied AI (EAI) system to make "good" decisions and prevent "bad" outcomes.
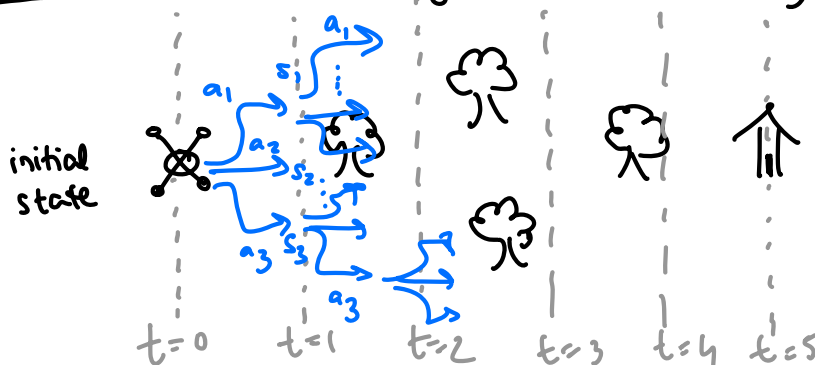


Our goal is to:

- mathematically <u>model</u> this decision-making process
- <u>quantify</u> the "goodness" of decisions
- <u>compute</u> these good decisions

This is where sequential decision-making (i.e. control theory) will provide us with a framework for answering these!

🔍 What makes sequential decision-making <u>hard</u>?

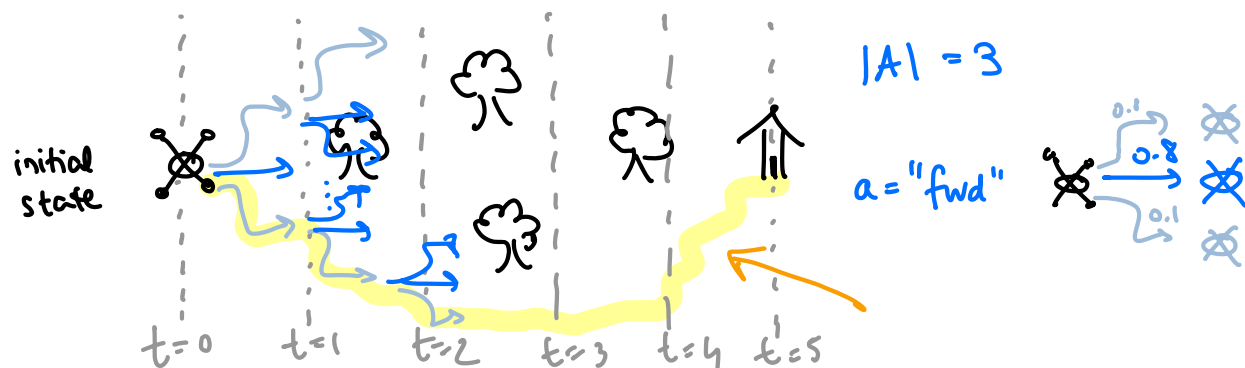📄 A1 naive solution grows exponentially in time horizon



$|A| = 3$ ⟨ up, straight, down

starting just from one state, you have $|A|^T$ possible seq. of decisions!

search: $\mathcal{O}(\left(|S| \times |A|\right)^T)$

$t=0 \quad t=1 \quad t=2 \quad t=3 \quad t=4 \quad t=5$

**A2** outcomes of taking actions can be stochastic (or unknown)



$|A| = 3$

$a = \text{"fwd"}$

initial state

t=0    t=1    t=2    t=3    t=4    t=5

0.1
0.8
0.1

**A3** Some aspects of the world are hard to **represent** & **predict**



cut

One framework that can help us describe these phenomena are:

**state space representations.**

state :  $x(t) \in \mathbb{R}^n$ ← continuous time   (often just $x$)  //  $s_t$ in the AI lit.
         $x_t \in \mathbb{R}^n$ ← discrete time

describes the minimal necessary characteristics of a system

v
$(x,y)$ $\overset{..)}{\bullet} \theta$      ex. position in x-y plane of 2D vehicle, orientation, speed

Control
Control / action :  $u(t) \in \mathbb{R}^m$  //  $a_t$ in the AI lit.
              AI    $u_t \in \mathbb{R}^m$

inputs that we choose @ each instance in time.

ex. joint torques, acceleration

output / observation : $y(t) \in \mathbb{R}^L$   // $o_t$ in the AI lit.
$\qquad\qquad\qquad\qquad y_t \in \mathbb{R}^L$

outputs that are actually measurable (typically through sensor)

ex. position through GPS, img. obs RGB camera,

! initially, we assume $y \equiv x$ "perfect observability" But it's not the case in general & it's a major safety challenge!

Control ↖

AI / RL ↗
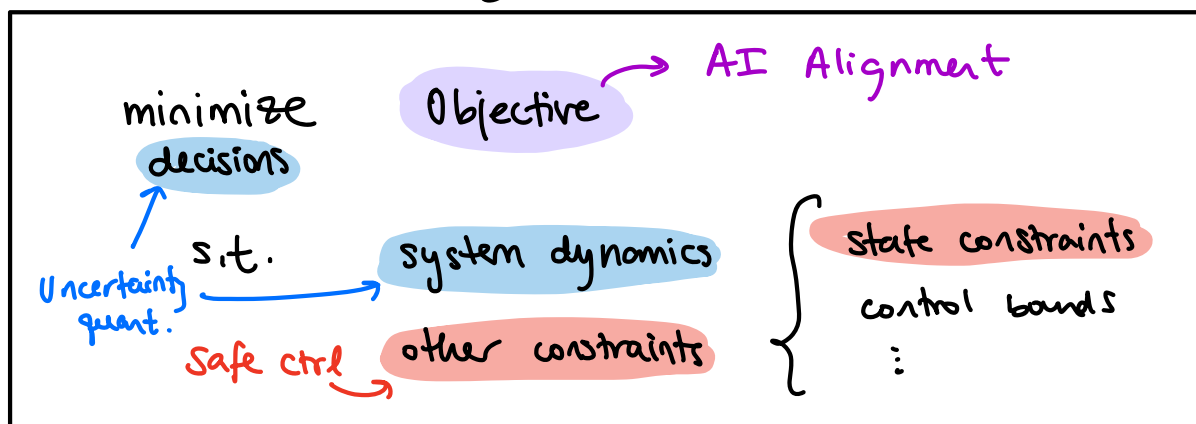
dynamics / transition : how the system evolves over time

| continuous-time $(t \in \mathbb{R})$ | discrete time $(t \in \mathbb{Z})$ |
|---|---|
| deterministic | next state, state, action |

deterministic   change in state ↘   state now ↙
$$\dot{x}(t) = f(x(t), u(t))$$
↑ action now

$$x_{t+1} = f(x_t, u_t)$$
(next state ↙, state ↙, action ↘)

stochastic
$$dX_t = f(X_t, u_t)dt + dW_t$$

diffusion models! ↗

Weiner Process (infintessimal Gaussian noise)

$$x_{t+1} \sim P(x_{t+1} \mid x_t, u_t)$$

## Optimal Decision - Making



minimize
decisions

Objective → AI Alignment

s.t. → system dynamics

other constraints

Uncertainty quant. ↑

Safe ctrl ↗

{ state constraints
control bounds
⋮

An optimization problem, but it has a temporal aspect to it and we want sequence of decisions that are optimal.

More formally, in this class we will see:

**Discrete-time:**

$\min\limits_{u_{0:T-1}} J(x_0, u_{0:T-1})$    $X \times U^T \to \mathbb{R}$    total cost of starting from $x_0$ & applying ctrl $u_{0:T-1}$

s.t. $x_{t+1} = f(x_t, u_t)$

$\underline{u} \leq u_t \leq \overline{u}$   ← control bounds   upper limit of ctrl.

lower limit   $\forall t \in \{0, \dots T-1\}$

**Continuous-time:**

$\min\limits_{u(\cdot) \in \mathcal{U}_0^T} J(x(0), u(\cdot))$

s.t. $\dot{x}(t) = f(x(t), u(t))$

$\underline{u} \leq u(t) \leq \overline{u}$

$\forall t \in [0, T]$

<u>ex.</u> of popular cost function (also called objective/reward)

$$J := \sum_{t=0}^{T-1} \underset{\text{running cost}}{L(x_t, u_t)} + \underset{\text{terminal cost}}{\ell(x_T)}$$

---

**BIG Q: How to solve?**

Each method listed below has its pros/cons -- BUT there are MANY solution strategies we can try! ☺

① Calculus of Variations   { converts constrained opt. → unconstrained via Lagrange multipliers

② Model Predictive Control (MPC)   { like ① but usually in discrete-time and it replan in "receeding horizon"

③ Dynamic Programming   { leverages the recursive structure in opt. control problems to compute policy

④ Reinforcement Learning   { same foundations as ③ but uses sim, func. approx, data to scale.
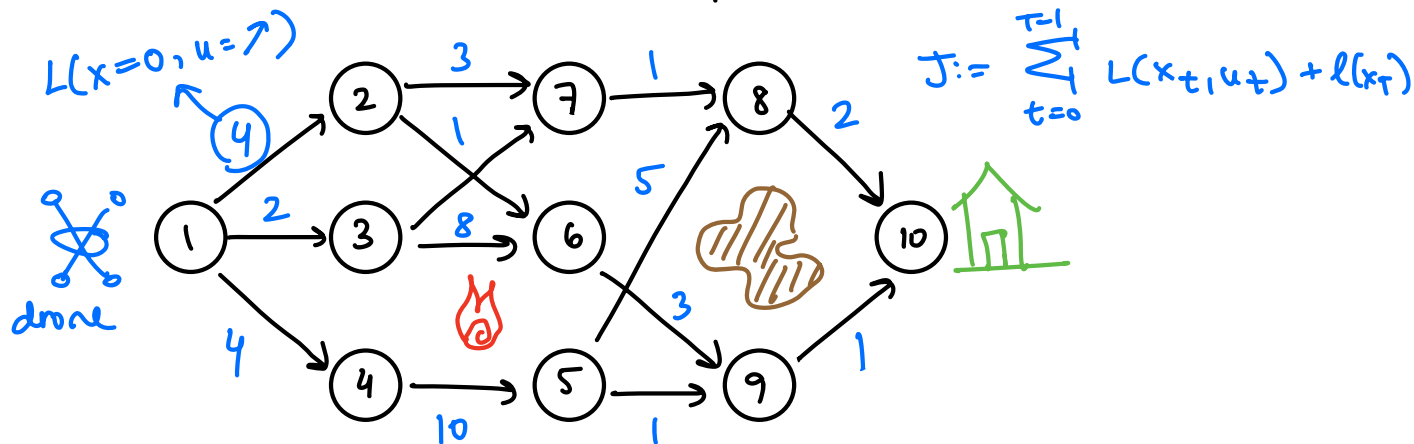
⇒ feature: unknown dynamics

# Dynamic Programming

We will study the method of dynamic programming to solve optimal control problems. Dynamic programming relies on the principle of optimality developed by Richard Bellman around 1958 when he was working @ the RAND corporation. Since then dynamic prog. has been used in CS, operations research, controls, robotics, and many other domains.
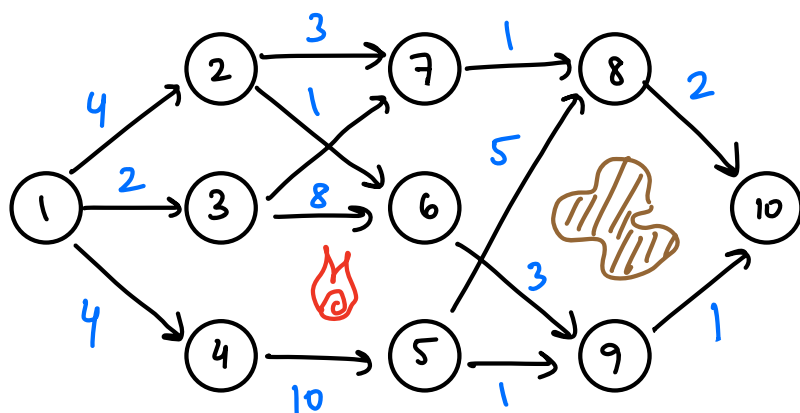
We can intuit dynamic programming via an example.

[ex.] DRONE rescue during California wildfires.

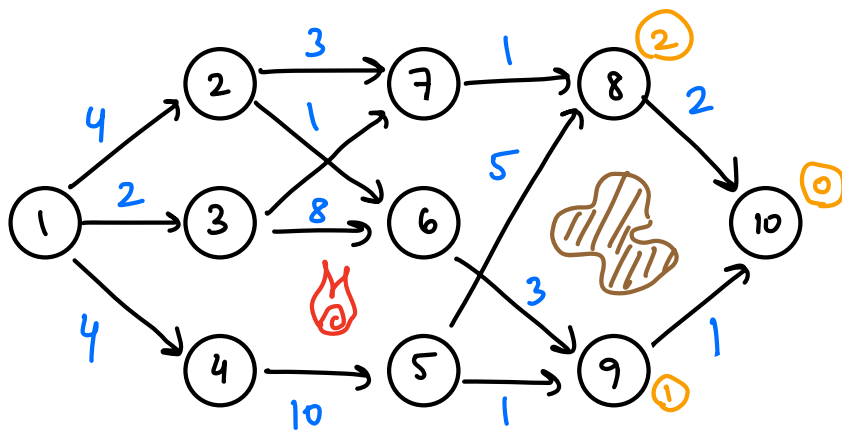⇒ need to find shortest path from ① → ⑩



$L(x=0, u=\nearrow)$

$$J := \sum_{t=0}^{T-1} L(x_t, u_t) + \ell(x_T)$$

drone

Q what's your strategy to solve this?

A work backwards!



First, if I'm already @ ⑩, then cost = 0

i.e. $\ell(x_T = 10) = 0$

Top diagram labels:
- Nodes: 1, 2, 3, 4, 5, 6, 7, 8, 9, 10
- Edges: 1→2: 4, 1→3: 2, 1→4: 4, 2→7: 3, 2→6: 1, 3→6: 8, 3→7 (5), 7→8: 1, 8→10: 2, 8→9: 3, 4→5: 10, 5→9: 1, 9→10: 1, 5→8: (via)
- Node costs: 8 has ②, 10 has ⓪, 9 has ①

Now, you look @ one timestep backwards @ node ⑧ and ⑨ and evaluate what is the cost to go from ⑧→⑩ & ⑨→⑩; pick the min!

Second diagram annotations:
- 3+3 = 6 / 1+4 = ⑤
- 1+2 = ③
- 2+0 = ②
- 1+3 = ④ / 8+4 = 12
- 3+1 = ④ 5
- 4+5 = 9 / 2+4 = ⑥ / 4+12 = 16
- 10+2 = 12
- 1+1 = ② / 5+2 = 7
- 1+0 = ①

It's more interesting for node ⑤. There are 2 ways to go ↗ ↘ but we only have to look one step ahead at the next nodes optimal cost to go!
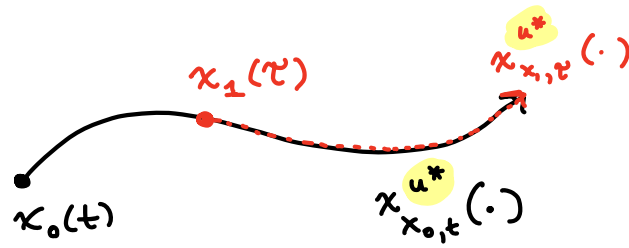
There are a few key properties of Dynamic programming

• DP gives you the optimal path from **all nodes** to node ⑩. You get intermediate sol^n "for free"

• Globally optimal solution

• DP gives computational gains over fwd. sim.

Let's understand underlying mathematical principle. DP relies on:

Principle of optimality: "In an optimal sequence of decisions or choices, each subsequence must also be optimal. Thus, if we take any state along the opt. state trajectory, then the remaining subtrajectory is also optimal"

In the example earlier, if we take any intermediate node along the optimal route, we still take optimal route to destination.

$$x_1(\tau)$$

$$x_{x,t}^{u^*}(\cdot)$$

$$x_0(t)$$

$$x_{x_0,t}^{u^*}(\cdot)$$

Let's write this principle down mathematically:

We want to solve:

"value function" $\leftarrow$ $V_t(x_t)$ $=$ $\min\limits_{u_{t:T-1}} J_t(x_t, u_{t:T-1})$ $= \sum\limits_{\tau=t}^{T-1} L(x_\tau, u_\tau) + \ell(x_T)$

define $V_t$ as storing the best-case "cost-to-go" from $x$ @ time $t$ to the end.

Let's expand out the RHS over time w/ our cost function:

depends on $(x_t, u_t)$

$$V_t(x_t) := \min\limits_{u_{t:T-1}} \left\{ L(x_t, u_t) + L(x_{t+1}, u_{t+1}) \cdots L_{T-1}(x_{T-1}, u_{T-1}) + \ell(x_T) \right\}$$

$$= \min\limits_{u_{t:T-1}} \left\{ L(x_t, u_t) + J_{t+1}(x_{t+1}, u_{t+1:T-1}) \right\} \quad \text{cost from next state } x_{t+1} \text{ onwards}$$

$$= \min\limits_{u_t} \left\{ L(x_t, u_t) \right\} + \left[ \min\limits_{u_{t+1:T-1}} J_{t+1}(x_{t+1}, u_{t+1:T-1}) \right] \quad \text{pull other min inside b/c } u_t \text{ only infl. } L \text{ \& } x_{t+1}$$

$:=V_{t+1}$    principle of optimality!

$$= \min\limits_{u_t} \left\{ L(x_t, u_t) + V_{t+1}(x_{t+1}) \right\}$$

just restrict ourselves to optimal trajectories starting from the next state

$$= \min\limits_{u_t} \left\{ L(x_t, u_t) + V_{t+1}(f(x_t, u_t)) \right\}$$

$u_t$ $\leftarrow$ no longer need to optimize over sequence, only current action!

w/ terminal condition $V_T(x) = \ell(x)$

Bellman Equation:

$$V_t(x_t) = \min_{u_t}\left[L(x_t, u_t) + V_{t+1}(x_{t+1})\right],$$

$$V_T(x_T) = \ell(x_T)$$

- The beauty is this lets us decompose decision-making problems into smaller subproblems and solve recursively, pointwise optim. over ctrl.

- $V(\cdot)$ is typically hard to solve in closed-form for most dynamical systems but for some you can!

  Exercise (offline): Linear Quadratic Regulator (LQR)

  i.e $x_{t+1} = A x_t + B u_t$ (lin. dyns.)

  $$L(x_t, u_t) = x^T Q x + u^T R u$$
  (quad. cost)