

Last Time

- sequential decision-making
- Dynamic programming

Lecture 2

EAIS SP'25

Andrea Bajcsy

This Time:

- what makes safe decision-making "hard"?
- safety filters

ANNOUNCEMENTS:

- HW #1 Released! It's on today's topic - safety filters ;)
- Ken's OH's:
Tuesdays, 1-2 pm, NSH 1511.

Attendance Question of the Day:

Why is dynamic programming "helpful" compared to forward simulation?

So far, we did a recap of mathematical modeling frameworks that let us think about decision-making.

$$\underset{u_{0:T}}{\text{maximize}} \quad J(x_0, u_{0:T}) \quad \leftarrow \text{objective}$$

$$\text{s.t.} \quad x_{t+1} = f(x_t, u_t) \quad \forall t \quad \leftarrow \text{dynamics}$$

$$u_t \in \mathcal{U} \quad \leftarrow \text{ctrl sounds}$$

$$x_t \notin \mathcal{F} \quad \leftarrow \text{state constraints}$$

This is what we will focus on in the first part of the class

$$x_t \notin \mathcal{F}$$

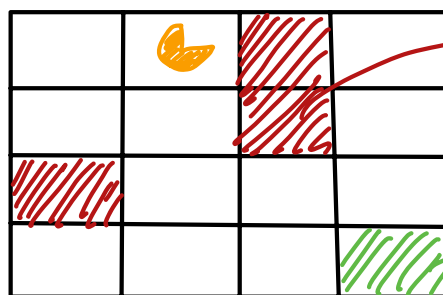
Q Give me examples of failure sets $\mathcal{F} \subset \mathcal{X}$ that might matter in:

- robot manipulator cleaning cluttered counter top
- drone flying through city

Before we go on, let's discuss why is ensuring that a robot makes decisions s.t. $x_t \notin \mathcal{F}$ hard?

I mean, if we know the failure set, isn't this enough?

For example, in Intro AI you'll see gridworld



← if going to hit \mathcal{F} : don't;

BUT, for EAI systems from this class, ensuring $x_t \notin \mathcal{F}$ is harder.

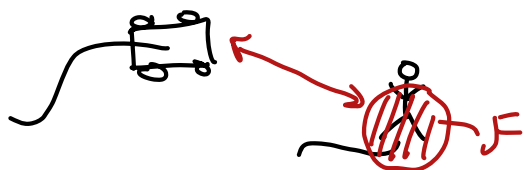
Reason #1. In practice, our decision-making "stack" or "policy" can be arbitrarily complex, or opaque.


e.g. E2E Diffusion Policy 

VLM - web agent.

Reason #2. other (strategic) agents

homicidal chauffeur problem (1971)



- Attacker or adversary 
- environment (ice, wind, ...)

Reason #3: Uncertainty

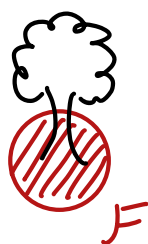
Even though we represent our system via a mathematical model, a model will never perfectly reflect reality.

Statistician George Box:

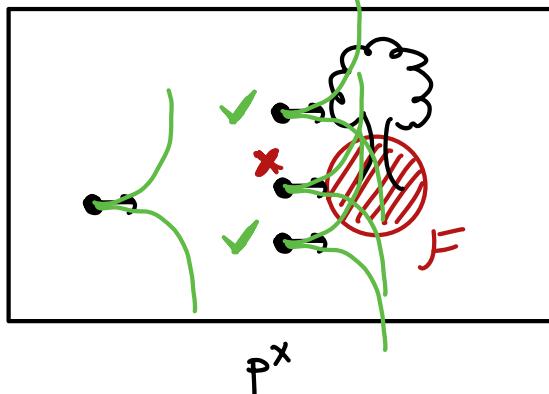
"All models are wrong, but some are useful"

Reason #4: Inevitable Failure

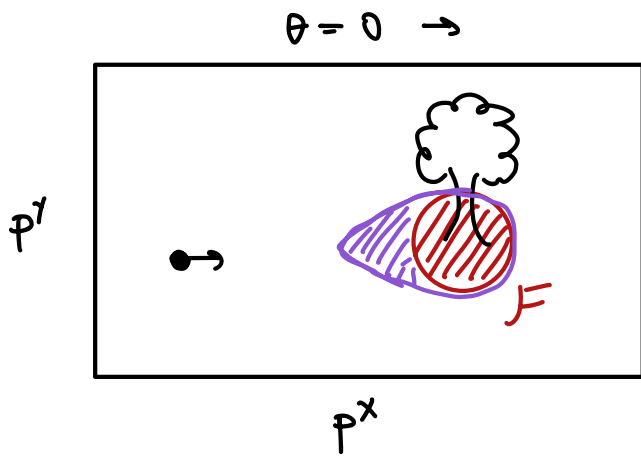

you biking home



p^y



Your brakes don't work! can only turn.



These purple states are the inevitable failure states

↓ also called unsafe set

How do we tackle Problems #1-4 rigorously but also practically?

IDEA: SAFETY FILTERS

Safety Filtering

Goal A "wrapper" we can put around ANY "base" policy / decision-making system

$$\pi \rightarrow u_{0:T}, \max_{u_{0:T}} J(x_0, u_{0:T})$$

(1) monitor if the system is @ "risk"

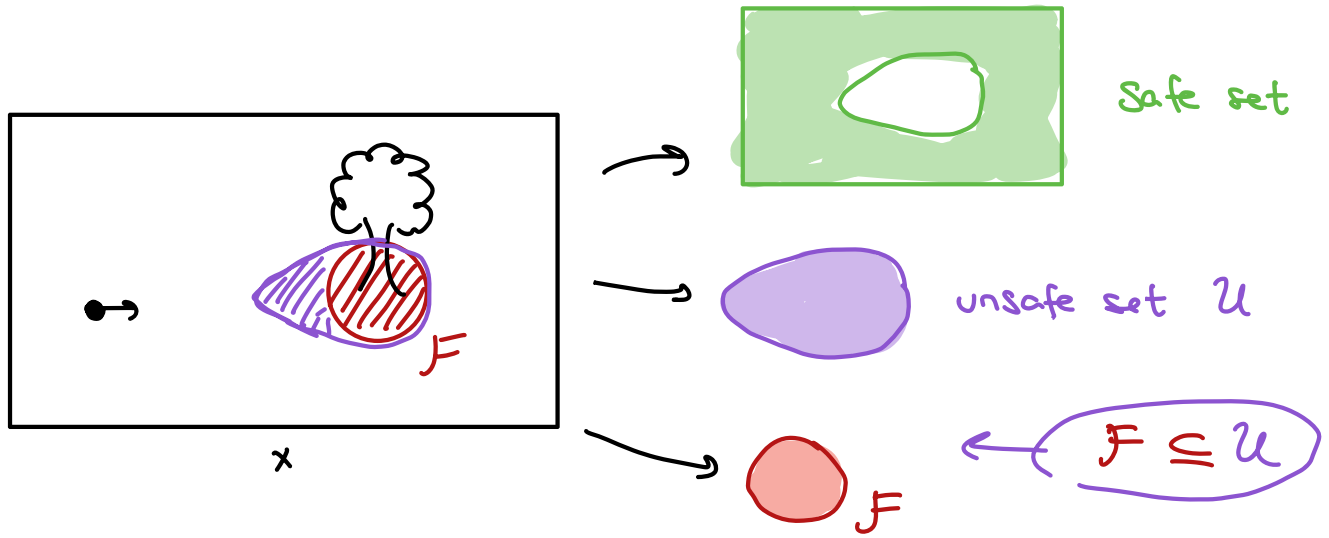
(2) adjust the robot's base policy @ runtime to prevent future failure. ($x_t \notin F \forall t \in \{0, 1, \dots, T\}$)

The idea is that we continue applying our nominal strategy for decision-making (e.g. VLM, Diff., MPC, ...) until safety is at risk; otherwise, apply a safety controller.

ex. SIMPLEST STRATEGY!

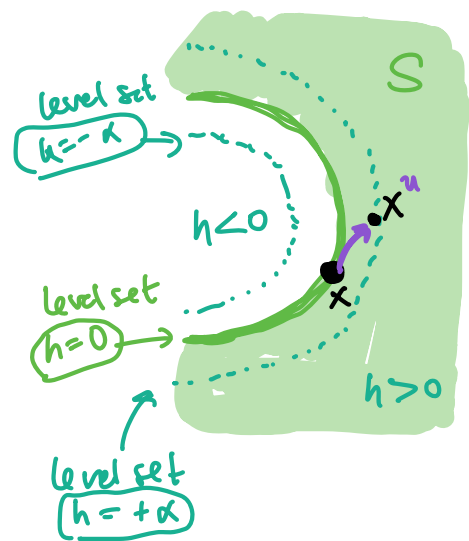
$$u^*(x) = \begin{cases} \pi^{\text{nom}}(x) & \text{if system is safe} \\ \pi^{\text{safe}}(x) & \text{if safety @ risk} \end{cases}$$

Q How do we know safety is @ risk, and the safe ctrl.?



! ASSUME we have a known safe set $S \subset X$ that :

- $S \cap F = \emptyset$ ↪ continuously differentiable fn $h: X \rightarrow \mathbb{R}$
- $$\begin{aligned} x \in S &\iff h(x) \geq 0 \\ x \in \partial S &\iff h(x) = 0 \end{aligned}$$
} encode the set via this func.
"boundary of S"



If we have access to such a set S and $h(x)$ function, then any policy $u := \pi(x)$ such that ↪ you are @ state x apply action u you get to state $x^u(t+\delta)$ in some small t -step δ .

$$h(x^u(t+\delta)) \geq 0$$

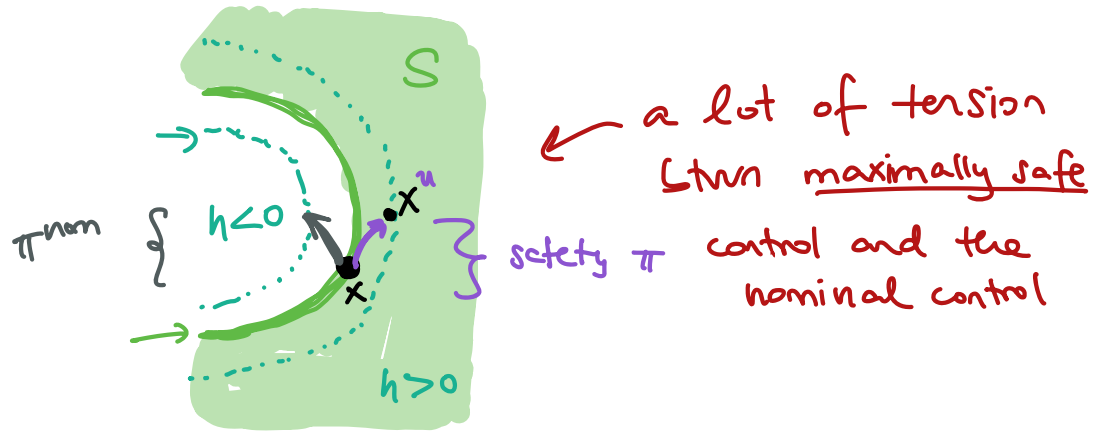
then this policy can prevent the system from leaving S (i.e. $\pi(x) \equiv \pi^{safe}$)

In this context, you may have heard of $h(x)$ function called a control barrier function (CBF) for S .

LEAST RESTRICTIVE SAFETY FILTER

$$u^*(x) = \begin{cases} \pi^{\text{nom}}(x) & \text{if } h(x) > 0 \text{ (i.e. } x \in S) \\ \pi^{\text{safe}}(x) & \text{if } h(x) = 0 \text{ (i.e. } x \in \partial S) \end{cases}$$

Q Problem: can lead to switch (aggressively) btwn. the nominal policy & safety policy



A key insight of CBF is an alternative safety filtering law which looks for similar control to the nominal one that is also safe.

$$u^*(x) = \arg \min_u \|u - \pi^{\text{nom}}(x)\|_2^2 \quad \left. \begin{array}{l} \text{"stay close to } \pi^{\text{nom}} \\ \text{while being safe"} \end{array} \right\}$$

s.t. u is safe

Q What is this? What is set of safe ctrl's?

A We know that

$$u \text{ is safe } (\Leftrightarrow) h(x^u(t+\delta)) \geq 0$$

from before.

$$h(x^u(t+\delta)) \geq 0$$

← Taylor series expansion

$$h(x^u(t+\delta)) \approx h(x(t)) + (t+\delta-t) \frac{dh(x(t))}{dt} \geq 0$$

$$= h(x(t)) + \delta \left[\frac{\partial h}{\partial x} \cdot \frac{\partial x}{\partial t} \right] \geq 0 \quad \leftarrow \text{chain rule}$$

$$= h(x(t)) + \delta \left[\frac{\partial h}{\partial x} \cdot f(x, u) \right] \geq 0$$

Now, we have the following safety filter:

$$\begin{aligned} u^*(x) = \arg \min_u & \| u - \pi^{\text{nom}}(x) \|_2^2 \\ \text{s.t. } & h(x(t)) + \delta \left[\frac{\partial h}{\partial x} \cdot f(x, u) \right] \geq 0 \end{aligned} \quad (*)$$

Some systems make solving fast! (i.e. a quadratic program)

CONTROL AFFINE:

$$f(x, u) := f_1(x) + f_2(x)u$$

This makes our optimization problem (*) a quadratic program!

→ objective is quadratic in u

→ constraint is linear in u

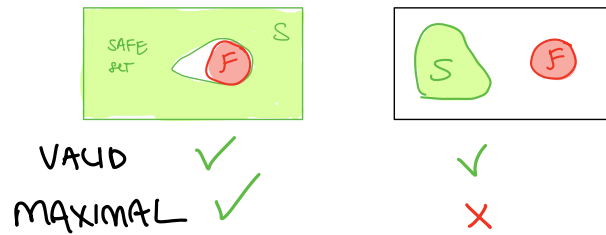
↓ solve fast online!

$$h(x(t)) + \delta \left[\frac{\partial h}{\partial x} f_1(x) + \frac{\partial h}{\partial x} f_2(x) u \right] \geq 0$$

BUT, let's talk about the ravids & elephant in the room

⚠️ Obtaining a VALID and not overly conservative safe set S is really challenging for general systems

and $h(\cdot)$ function!



⚠️ The CBF framework we have seen today doesn't handle uncertainty/disturbances robustly, and we also haven't talked about control input constraints. Disturbances typically break the assurance of S we have seen so far.

We will tackle these challenges head-on during the next lecture:

synthesizing (i.e. computing) safe sets & safety filters

↳ we will talk about a general, computational framework for getting S and π^{safe} and $h(\cdot)$ that

- ① is guaranteed to be VALID & MAXIMAL
- ② naturally handles disturbances robustly
- ③ is compatible w/ modern comp. tools

$\left. \begin{array}{l} \text{RL} \\ \text{SSL} \end{array} \right\} !$