

Last Time

- computing safety filters
- HJ reachability

Lecture 4

EAS SP'25

Andrea Bajcsy

This Time:

- Recap + hands-on-exercise
- code demo!
- robust safety

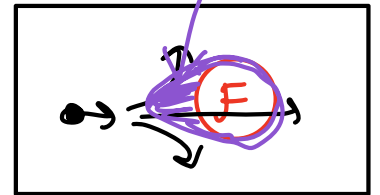
Recap - HS Reachability for Computing Safety Filters

Last time, we formalized the problem of safety satisfaction as an optimal control problem:

Continuous-time Safety Problem

$$V(x, t) := \max_{u(\cdot) \in \mathcal{U}_t^T} \min_{\tau \in [t, T]} l(x(\tau))$$

signed dist. func. s.t.
 $F = \{x : l(x) < 0\}$
 $l(x(\tau))$



"Safety Bellman Backup" (i.e. HJ-VI)

"remember if failing now"

$$\min \left\{ \overbrace{l(x) - V(x, t)}^{\text{"remember if failing now"}}, \underbrace{\frac{\partial V}{\partial t} + \max_{u \in \mathcal{U}} \frac{\partial V}{\partial x} \cdot f(x, u)}_{\text{"try best not to fail in future"}} \right\} = 0$$

$$V(x, T) = l(x)$$

"try best not to fail in future"

Discrete-time Safety Problem

$$V_t(x_t) := \max_{u_t, u_{t+1}, \dots, u_T} \min_{\tau \in \{t, t+1, \dots, T\}} l(x_\tau)$$



"Safety Bellman Backup" (discrete-time)

remember failing now"

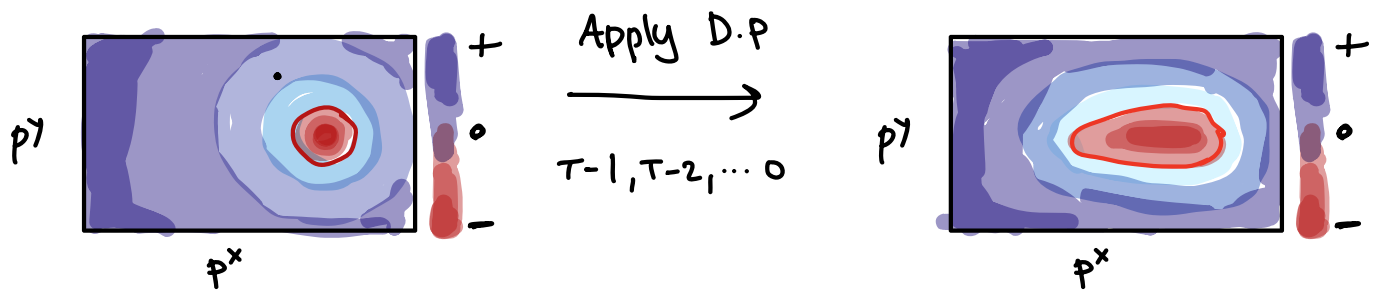
$$V_t(x_t) = \min \left\{ \overbrace{l(x_t)}^{\text{"remember failing now"}}, \max_{u_t \in \mathcal{U}} V_{t+1}(f(x_t, u_t)) \right\}$$

$$V_T(x_T) = l(x_T)$$

"try best not to fail in future"

$$V(x, T) = \ell(x) \downarrow$$

$$V(x, 0)$$



Getting our safety monitor (i.e. the unsafe set boundary)

By design, the zero-sublevel set of $V(\cdot, \cdot)$ encodes our unsafe set!

$$\text{Unsafe set} \equiv \text{BRT}(t) = \{x: V(x, t) < 0\}$$

Getting safety-preserving policy

We can also compute the optimal safety policy via the optimal value function too!

cont. time \downarrow

$$u^{\text{safe}}(x, t) = \arg \max_{u \in \mathcal{U}} \frac{\partial V}{\partial x} \cdot f(x, u)$$

discr. time \downarrow

$$u_t^{\text{safe}}(x_t) = \arg \max_{u_t \in \mathcal{U}} V_{t+1}(f(x_t, u_t))$$

Safety Filter

$$u(x) = \begin{cases} \pi^{\text{nom}} & \text{if } V > 0 \\ u^{\text{safe}} & \text{if } V \approx 0 \end{cases}$$

Exercise: Compute a Backwards Reachable Tube (i.e. unsafe set)

We will operate in discrete state + time to build intuition:

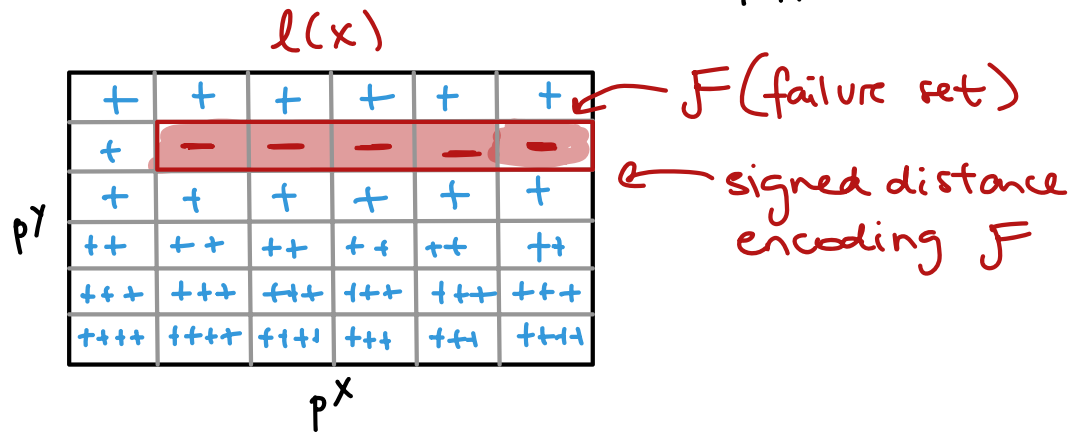
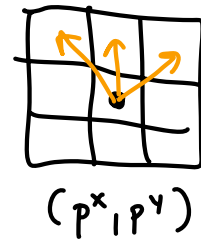
"Safety Bellman" Eqn:

$$V_t(x) = \min \{ l(x), \max_{u \in \mathcal{U}} V_{t+1}(f(x, u)) \}$$

$$V_T(x) = l(x)$$

• system state: $x = \begin{bmatrix} p^x \\ p^y \end{bmatrix} \in \mathcal{X}$

• system ctrl: $u \in \mathcal{U} = \{ \leftarrow, \uparrow, \rightarrow \}$



We initialize $V_T(x) = l(x)$. Compute $V_{T-1}(x)$ and $V_{T-2}(x)$.

Solution:

+	+	+	+	+	+
+	-	-	-	-	-
+	+	+	+	+	+
++	++	++	++	++	++
+++	+++	+++	+++	+++	+++
++++	++++	++++	++++	++++	++++

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

+	+	+	+	+	+
+	-	-	-	-	-
+	+	-	-	-	-
++	++	+	+	+	+

$$V_{T-1}(x_k) = \min\{-, +\} = -$$

$$V_{T-1}(x_{T-1})$$

$$V_{T-1}(x_j) = \min\{+, +\} = +$$

$$V_{T-1}(x_i) = \min\{+, \max V_T(f(x_i, u))\} = -$$

		-	-	-	-
		-	-	-	-
		-	-	-	-
		-	-	-	-

$$V_{T-2}(x_{T-2})$$

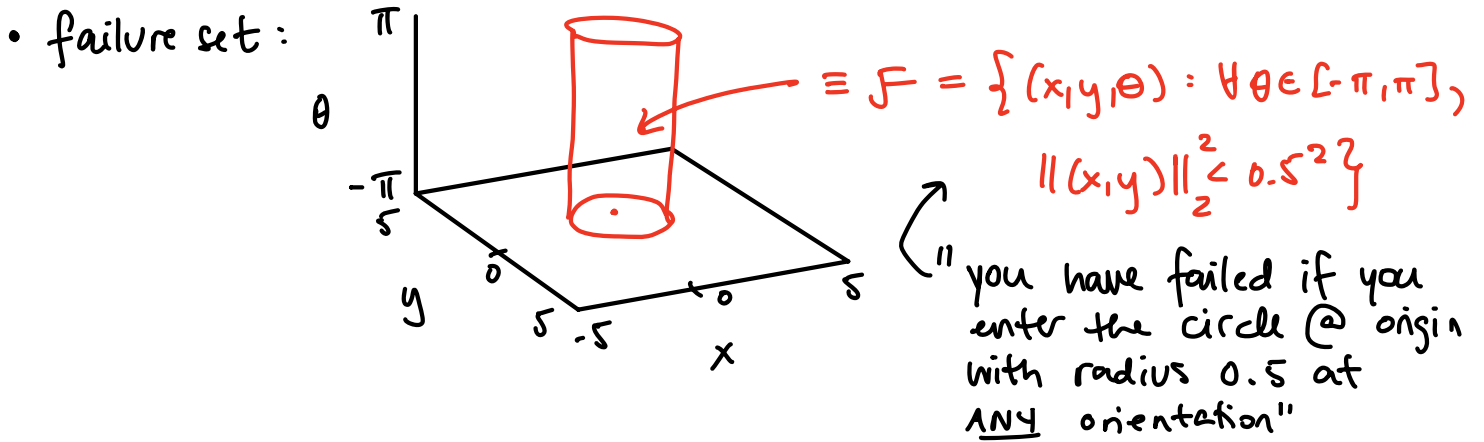
$$V_{T-2}(x_m) = \min\{++, -\} = -$$

Code Demo (in MATLAB solver helperOC + LevelSetToolbox)

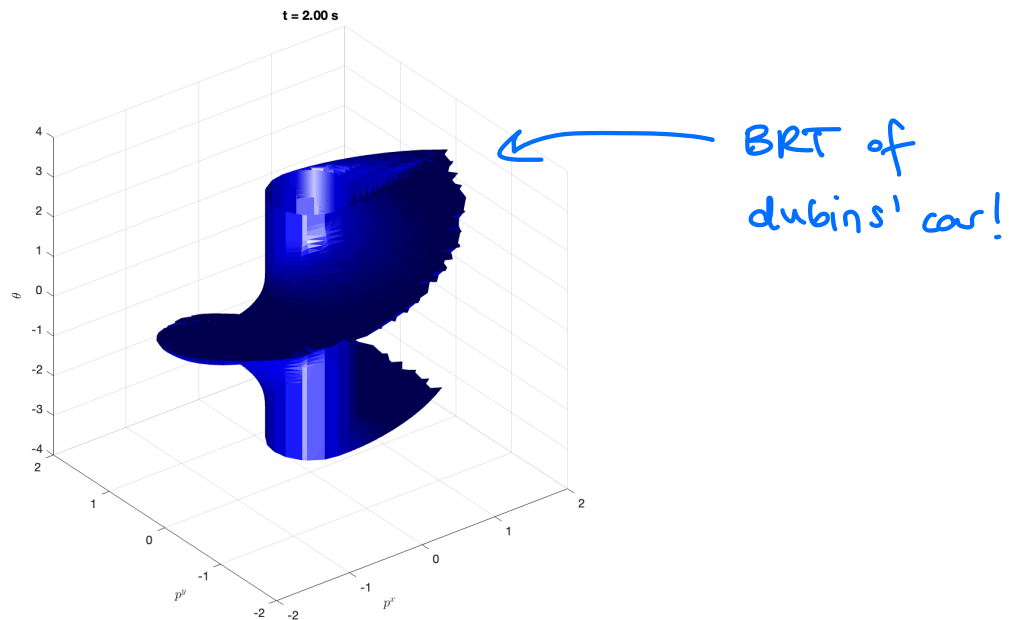
• system:

$$\begin{aligned} \dot{x} &= v \cos \theta = 1.5 \text{ m/s} \\ \dot{y} &= v \sin \theta \\ \dot{\theta} &= u \rightarrow u \in [-0.5, 0.5] \end{aligned}$$

} Dubins' car



• computed unsafe set (i.e. backwards reachable tube):



Ⓛ Way easier as an engineer to specify failure set F than to specify BRT. That's why we want to compute BRT given F !

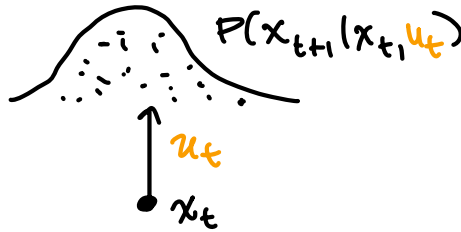
Robustifying Safety

So far, we have assumed that our dynamical system perfectly evolves via $\dot{x} = f(x, u)$ with no uncertainty.

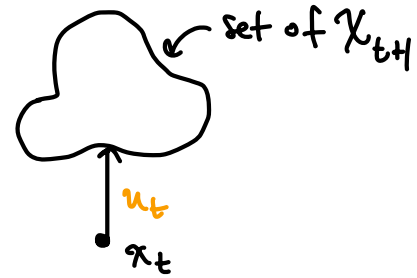
This isn't realistic for many real-world scenarios (e.g. friction!)



No uncertainty



Probabilistic



Non-deterministic

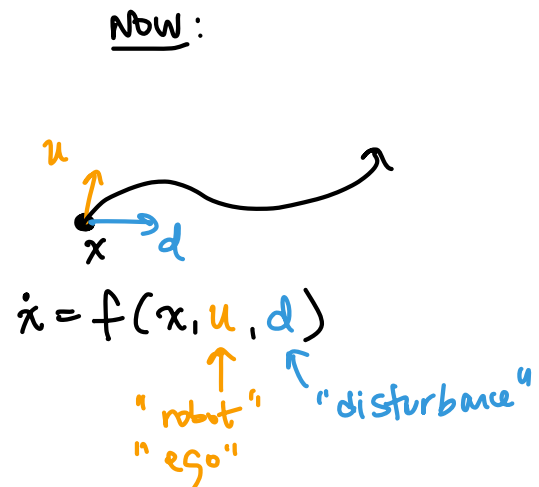
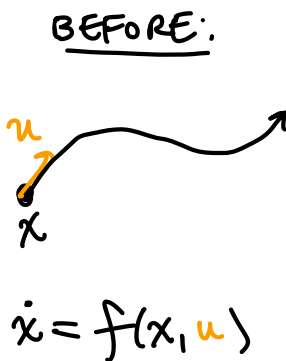
There are two ways to model uncertainty:

1) probabilistic uncertainty (i.e. "I have observed data")

2) non-deterministic uncertainty (i.e. "I have minimal additional info.")

How should we handle the design of our safety filter (+ analysis) to handle uncertainty?

Typically, we do this via modelling another "input" that influences the state evolution:



In robust safety, we take a non-deterministic view of uncertainty and assume that $d \in \mathcal{D}$ is chosen from some bounded set and we want our robot to be ROBUST to the WORST POSSIBLE sequence of d 's!

ROBUST BACKWARDS REACHABLE TUBE (BRT) of a set $F \subset X$

and dynamical system $\dot{x} = f(x, u, d)$ is:

$$\text{BRT}(t) := \left\{ x \in X : \underbrace{\forall u(\cdot) \in \mathcal{U}_t^T}_{\text{for all things the robot could do}}, \underbrace{\exists d(\cdot) \in \mathcal{D}_t^T}_{\text{there is something the disturbance could do ...}} \right. \\
 \left. \underbrace{x_{\tau}^{u,d}(\tau) \in F}_{\text{that leads the robot to failure}} \text{ for some time } \tau \in [t, T] \right\}$$

This is the set of all starting states from which no matter the controller's effort, the disturbance can push system into F . The way we will formulate an "optimal control" problem whose solution represents this unsafe set will be via:

<u>zero-sum</u> there is a winner + loser	<u>dynamic</u> game <u>evolves</u> over time	<u>games</u> result/outcome depends on 2+ players/inputs
--	--	--