Last Time

☐ dynamic games

This Time:

☐ compare BRT with + w/o disturbance

☐ scaling computation : RL

☐ scaling computation: SSL

Announcement    HW #1   due  today !

# Numerical Comparison of Robust vs. Non-Robust Unsafe Set

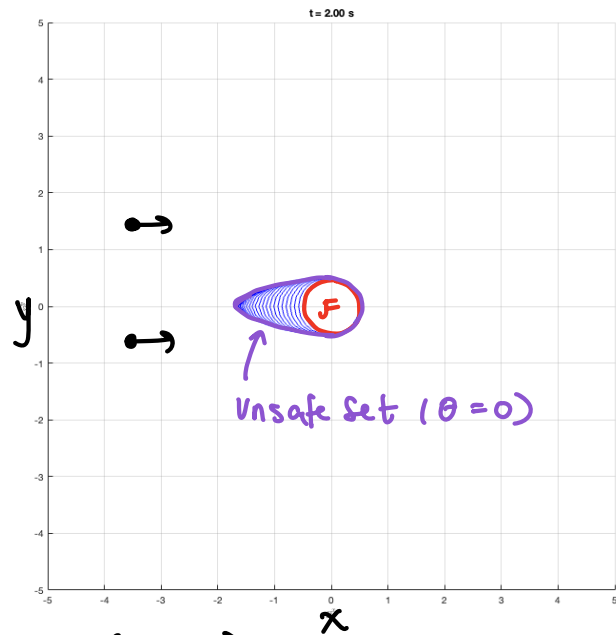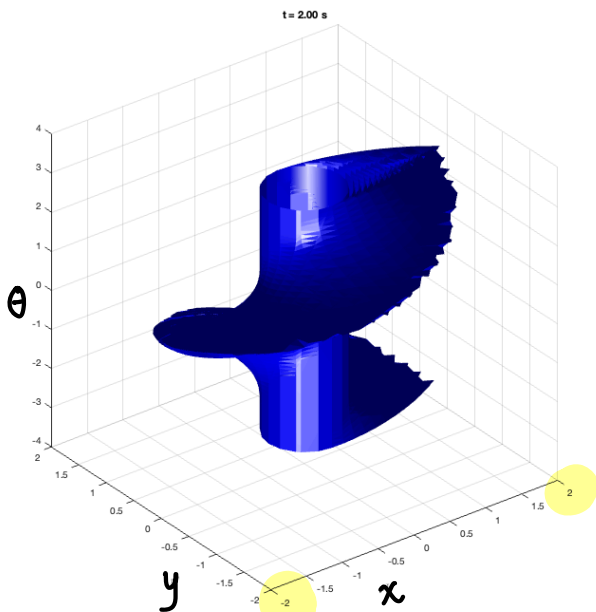Recall: state is $(x, y, \theta)$ of Dubins' car

dynamics are
$$\dot{x} = v \cos \theta + d_x \in [-0.5, 0.5]$$
$$\dot{y} = v \sin \theta + d_y \equiv D$$
$$\dot{\theta} = u \in [-0.5, 0.5] \equiv \mathcal{U}$$

with $v = 1.5$

failure is cylinder @ origin with radius 0.3

Non-Robust BRT

$$\max_{u} \min_{\tau \in [t, \tau]} \ell(x(\tau))$$



t = 2.00 s



t = 2.00 s
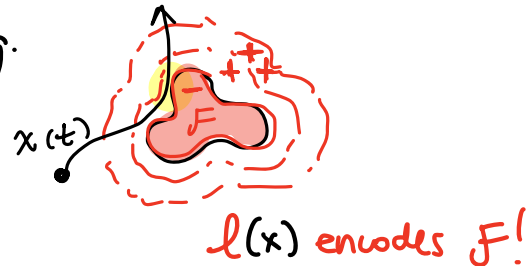
Unsafe set ($\theta = 0$)

Robust BRT

$$\max_{\pi_u} \min_{\pi_d} \min_{\tau \in [t, \tau]} \ell(x(\tau))$$



t = 2.00 s



t = 2.00 s

Unsafe set ($\theta = 0$)

_RECAP:_

So far, we have formalized what is a safe set + safe controller. We also transformed this into an optimal control (or dynamic game for robustness) problem that we can solve via dynammic programming.

$x(t)$

$\ell(x)$ encodes $F$!

Now that we have the foundations, we can talk about _practical challenges_ and _frontiers_ of decision-theoretic safety!

PRACTICAL CHALLENGES + RELATED RESEARCH FRONTIERS

1) scaling the computation of BRT's (i.e. value function!)

"FRONTIERS II"

2) increasing the flexibility of the safe/unsafe set (e.g. representation is parameterized, use data to inform BRT)

3) specifying more complex failure representations (e.g. "don't go through caution tape", "don't spill")

4) break the perfect state assumption

5) analyze more complex dynamical systems

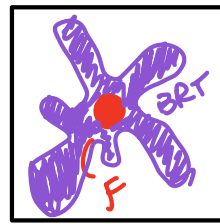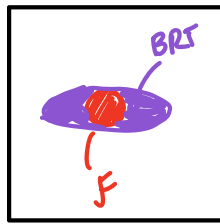[TODAY] we will focus on (1) — scaling the computation!

At the highest level, the key approach to scaling we will study is:

NEURAL APPROXIMATIONS to the value function

These "neural reachable tubes" will bake safety into the training process (which will leverage data) via careful construction of the signal we use for learning the neural approximation.

NN's are "agnostic" to grid resolution, and so the memory + the learning complexity scales with the underlying complexity of the set (not the resolution)



simple underlying BRT shape! ⟹ easy to represent + learn for NN!

BRT

F

BRT

F

← harder underlying BRT shape ⟹ more complex to represent + learn!

Broadly, there are two "paradigms" for solving for our neural reachable tubes:

- reinforcement learning (RL)
- self-supervised learning (SSL)

## Reinforcement Learning

RL approaches leverage advances in high-fidelity simulators that model hard-to-write-down $\dot{x} = f(x,u)$ systems and they use "rollouts" (i.e. data of the robot executing its behavior in the (simulated) environment to approximate an optimal control / sequential decision-making problem.

! KEY CHALLENGE: RL algorithms are designed for discounted sum-of-reward problems. BUT, safety isn't about failing "a little" on average; its about preventing failure AT ALL.

Typically, the RL problem is posed as approximating

$$V(x) := \max_{\pi} \sum_{t=0}^{\infty} \gamma^t r(x_t, u_t)$$

determines importance of future rewards: $\gamma = 0 \Rightarrow$ myopic
$\gamma \to 1 \Rightarrow$ optimal

← discount factor $\gamma \in [0,1)$

The corresponding Bellman backup for this problem (w/ discrete-time determistic dynamics) is:

$$ \circledast \qquad V(x) = \max_{u \in \mathcal{U}} r(x, u) + \gamma V(f(x, u)) $$

An important property that allows many RL algorithms to work is the fact that $\circledast$ is a <u>contraction mapping</u>. Intuitively, this means that successive applications of $\circledast$ will <u>converge</u> to a unique <u>fixed point</u>/solution.

$$ \nearrow V^*(x) $$

The problem is that our (infinite-horizon) safety problem <u>doesn't</u> does not induce a contraction mapping in the Bellman backup.

$$ V(x) := \max_{\pi} \min_{t \geq 0} \ell(x_t^{\top}) \qquad \leftarrow \text{o.c. problem} $$

$$ \Downarrow $$

$$ \leftarrow \text{Bellman Backup} $$

$$ (\circledast\circledast) \qquad V(x) = \min \left\{ \ell(x), \max_{u \in \mathcal{U}} V(f(x, u)) \right\} $$

How do we discount $(\circledast\circledast)$ which has a min over time....

<u>ONE IDEA</u>: Fisac*, Lugovoy*, Rubies-Royo, Gosh, Tomlin. ICRA 2019.

$$ \circledast \qquad V(x) = (1 - \gamma) \ell(x) + \gamma \min \left\{ \ell(x), \max_{u \in \mathcal{U}} V(f(x, u)) \right\} $$

Intuitively, they interpret $\gamma \in [0, 1)$ as the probability of the episode continuing. If it continues, then the value will be the RHS $\circledast$ "minimum of the future" with probability $(1 - \gamma)$, it may end (e.g. transition into a terminal state) and so your current value is $\ell(x)$.

Let's prove that $\circledast$ is contraction mapping!

**Thm** The backup operator $B[V]$:

$$B[V] := (1-\gamma)\ell(x) + \gamma \min\left\{\ell(x), \max_{u \in U} V(f(x,u))\right\}$$

is a contraction mapping. Let $V, \tilde{V}: X \to \mathbb{R}$

There exists $K \in [0,1)$ s.t.

$$\|B[V] - B[\tilde{V}]\|_\infty \leq K\|V - \tilde{V}\|_\infty$$

$$\|x\|_\infty := \max_i |x_i|$$

**Proof:** Consider $\forall x \in X$

$$|B[V] - B[\tilde{V}]| =$$

next state $x' = f(x,u)$

$$\left|(1-\gamma)\ell(x) + \gamma \min\left\{\ell(x), \max_{u \in U} V(x')\right\} - \right.$$

cancel out

$$\left. \left[(1-\gamma)\ell(x) + \gamma \min\left\{\ell(x), \max_{u \in U} \tilde{V}(x')\right\}\right]\right|$$

pull out of abs. value b/c $\gamma \in [0,1)$

$$= \gamma\left|\min\left\{\ell(x), \max_{u \in U} V(x')\right\} - \min\left\{\ell(x), \max_{u \in U} \tilde{V}(x')\right\}\right|$$

$\leq \max_{u \in U} V(x')$ b/c of min with $\ell(x)$!

$\leq \max_{u \in U} \tilde{V}(x')$

$$\leq \gamma\left|\max_{u \in U} V(x') - \max_{u \in U} \tilde{V}(x')\right|$$

Without loss of generality, suppose the first max is larger:

$$|B[V] - B[\tilde{V}]| \leq \gamma\left|\max_{u \in U} V(x') - \max_{u \in U} \tilde{V}(x')\right|$$

**Lemma:**

$$\left|\max_a f(a) - \max_a g(a)\right|$$

$$\leq \gamma \max_{u \in U}\left|V(x') - \tilde{V}(x')\right|$$

$K \equiv \gamma$

$$\leq \max_a |f(a) - g(a)| \leq \gamma \sup_{x' \in X}\left|V(x') - \tilde{V}(x')\right| = \gamma\|V - \tilde{V}\|_\infty$$

look @ all next states & find the max

By def$^n$ of $\|\cdot\|_\infty$

Now that we have a contraction mapping, we can unlock RL algorithms, but with our safety-informed backup!

e.g. Safety Q-learning

$$Q_\theta(x, u) \leftarrow (1 - \alpha)Q_\theta(x, u) +$$

learning rate $0 < \alpha \leq 1$

$$\alpha \left[ (1-\gamma)\ell(x) + \gamma \min \left\{ \ell(x), \max_{u' \in U} Q_\theta(x', u') \right\} \right]$$

"weighted average"

↳ new value estimated from experience

$Q_\theta(x, u)$ — state-action value func. parameterized by $\theta$

---

## Self-Supervised Learning
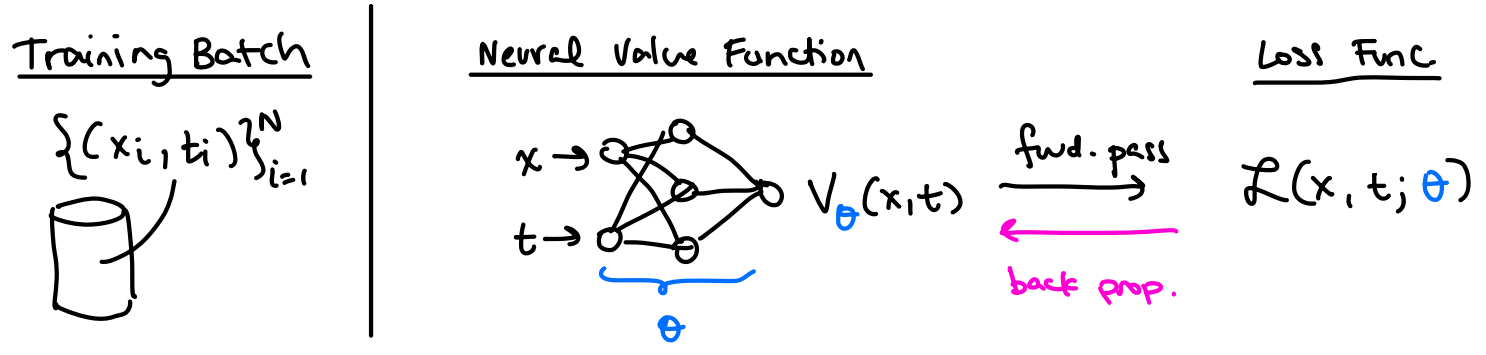
Still use neural reachable tubes, but pose SSL problem

IDEA. Bansal & Tomlin. "DeepReach". ICRA 2021.

¿ lets use the (continuous-time) HJB-VI equation as supervision!

If we found a good value function approx. $V_\theta(x, t)$, $\forall (x, t) \in \mathcal{D}$

$$\min \left\{ \ell(x) - V_\theta(x, t), \frac{\partial V_\theta}{\partial t} + \max_{u \in U} \nabla_x V_\theta(x, t)^T f(x, u) \right\} = 0$$

It must be that this PDE-like equation is equal zero!

---

| Training Batch | Neural Value Function | Loss Func |
|---|---|---|
| $\{(x_i, t_i)\}_{i=1}^N$ |  $x \to$ $t \to$ $V_\theta(x, t)$ $\xrightarrow{\text{fwd. pass}}$ $\xleftarrow{\text{back prop.}}$ $\theta$ | $\mathcal{L}(x, t; \theta)$ |

$$\mathcal{L}(x_i, t_i; \theta) = \overbrace{\text{HJ-VI Violation Error}}^{\text{the HJ-VI should} \equiv 0!} + \lambda \overbrace{\text{Initial Condition}}^{V(x,T) = \ell(x)}$$

$$= \mathcal{L}_1(x_i, t_i; \theta) + \lambda \, \mathcal{L}_2(x_i, t_i; \theta)$$

$$\mathcal{L}_1 = \sum_i \left\| \min\left\{ \ell(x_i) - V_\theta(x_i, t_i), \frac{\partial V_\theta(x_i, t_i)}{\partial t} + \max_u \nabla_x V_\theta(x_i, t_i)^\top f(x_i, u) \right\} \right\|$$

we know this $\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad}$ should $= 0$ when $V$ is optimal $\Rightarrow$ bigger the norm, the higher the loss!

$$\mathcal{L}_2 = \sum_i \| V_\theta(x_i, t_i) - \ell(x_i) \| \, \mathbb{1}\{t_i = T\} \quad \leftarrow \text{ground truth value only available @ final time.}$$