

Last time :

- implementing uncertainty quant.

Lecture 14

SP '26

Andrea Bajcsy

This Time:

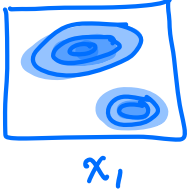
- system-level failures

Component vs. System-level Failures

So far, we have discussed a suite of UQ techniques for a data-driven model:

→ Density estimators

$\tilde{p}(x|\mathcal{D}): x_2$



x_1

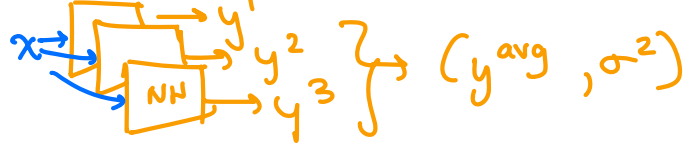
→ Cosine similarity:

$$s(x; \mathcal{D}) = \max_{\bar{x} \in \mathcal{D}} \frac{\bar{x}^T x}{\|\bar{x}\| \|x\|}$$

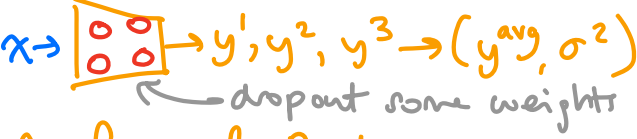
"gets most similar prior experience from \mathcal{D} to construct score"

inputs $x \rightarrow$ NN_{θ} \rightarrow predictions y


→ Ensembles



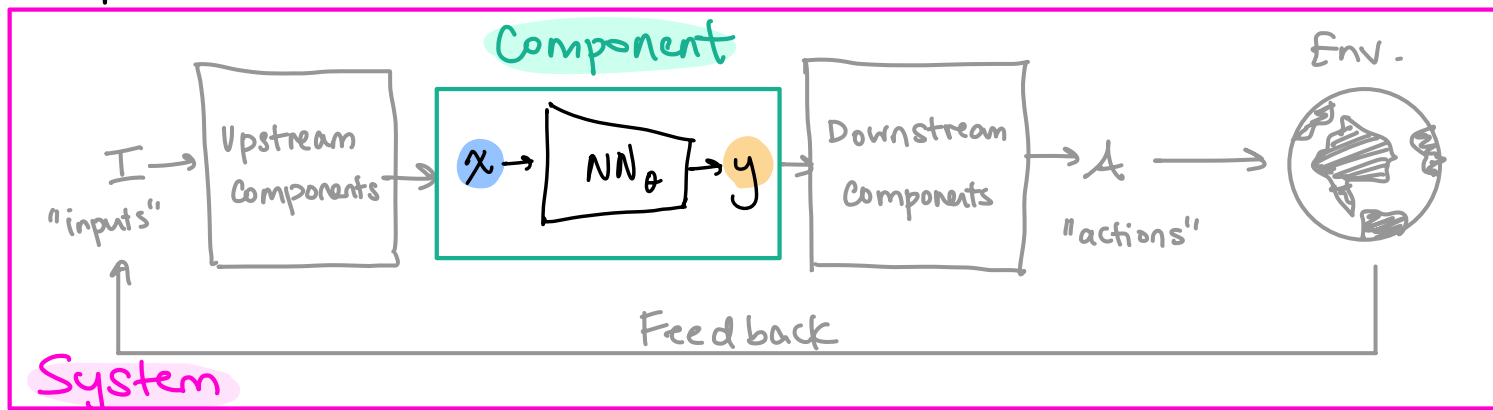
→ MC Dropout



→ Conformal Pred.



But, all these models live within a complex "autonomy pipeline" — ex. the NN could be perception module, trajectory forecaster, LLM task planner...



Q How can we model, identify, & repair component-level errors that cause system-level degradation?

MODEL

How do we mathematically model component-level errors that cause system-level degradation?

Broadly, all methods that study this problem end up tackling something that looks like this:

trajectory resulting from executing system when component sees input x

find s.t. $x \in X$ $\text{Cost} \left[\zeta^{NN(x)} \right] \geq \delta$

inputs into the learned component

function that evaluates performance / safety / risk / ... of closed-loop behavior.

Intuitively, we want to find inputs into our component that cause our system to degrade. Then we can define:

$$X^{\text{anomaly}} := \left\{ x : \text{Cost} \left[\zeta^{NN(x)} \right] \geq \delta \right\}$$

as the set of all "system-level anomalies".

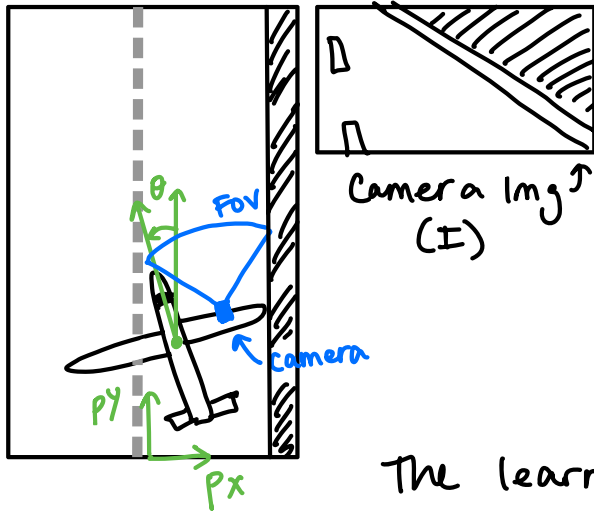
IDENTIFY

OK, but how do we practically instantiate this? Seems really hard (and it is!) but let's look @ two case studies to understand:

Case Study #1: Perception Module.

↳ From: Chakraborty & Bansal. "Discovering Closed-loop Failures of Vision-Based Controllers..." RA-L, 2023.

The setup: you have an aircraft w/ camera on right wing.

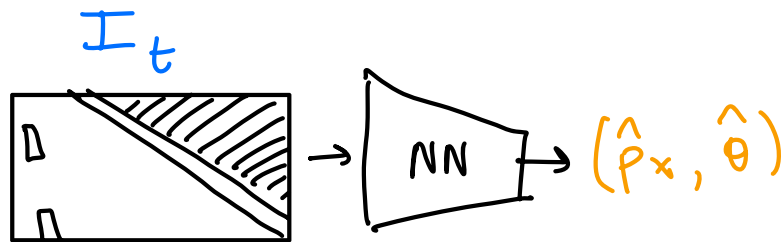


The true state + dyns. of

system are modelled as:

$$\begin{aligned} \dot{p}_x &= v \cos \theta & p_x &= \text{"cross track error"} \\ \dot{p}_y &= v \sin \theta & p_y &= \text{"downtrack error"} \\ \dot{\theta} &= u \leftarrow \text{control ang. vel.} & \theta &= \text{"heading error"} \end{aligned}$$

The learned perception model (is a CNN) takes as input the image I_t and returns an estimated $(\hat{p}_x, \hat{\theta})$



Then, a proportional controller (P-ctrl) steers the aircraft towards the centerline:

$$\pi(I_t) := \tan(-0.74 \hat{p}_x - 0.44 \hat{\theta}) = u$$

The failure states of the overall system are the aircraft leaving the runway: $\mathcal{F} = \{x: |p_x| \geq 10\}$

Thus, our goal is: find s.t. $\exists t \in [0, T], \{ \pi(I_t) \}_{t \in [0, T]} \in \mathcal{F}$

Ⓚ What type of problem does this remind you of?

[A] Key IDEA: we can pose this as a reachability problem!

But, searching over space of all $\mathcal{I}_{0:T}$ is too high-D...

We can exploit a trick here to reduce dimensionality:

In simulation, we know the sensor mapping from any

true state x_t to the image \mathcal{I}_t ! let $\mathcal{I}_t = S(x_t; E)$

be this sensor map from true state x_t in environment E .

So, we can rewrite:

$$u = \pi(\mathcal{I}_t) = \pi(S(x_t, E)) \Rightarrow \hat{\pi}(x_t)$$

This is the closed-loop, state feedback policy

Now, compute the BRT for all true states:

$$\text{BRT}_x := \{ x : \exists \tau \in [0, T] \ \gamma_x^{\hat{\pi}}(\tau) \in \mathcal{F} \}$$

And then transform it into the image space to get the set of all failure inputs:

$$\text{BRT}_{\mathcal{I}} = \{ \mathcal{I} : \mathcal{I} = S(x; E), x \in \text{BRT}_x \}$$

- Recipe:
- ① get closed-loop dyns. $\hat{\pi}$ via uniform sampling over state space \mathcal{X} & then rendering \mathcal{I} + getting u
 - ② computing $\text{BRT}_{\mathcal{I}}$ to get set of images that will lead to failures.

⚠ Compared to fwd-simulation methods:

→ reachability: 6.5 hours (most time on rendering)

→ fwd sim: ~ 67.5 (DAYS)

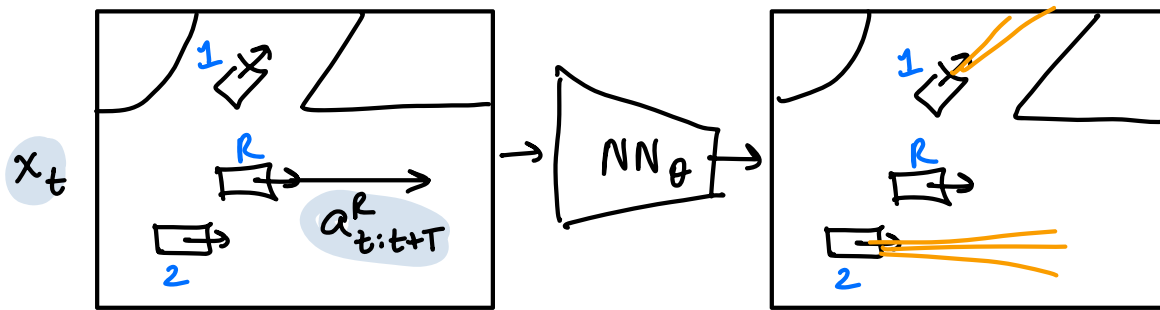
CASE STUDY #2: Trajectory Forecasting

↳ From: Nakamura, Tian, Bajcsy. "Not All Errors Are Made Equal." CoRL, 2024.

The Setup: We have a learned trajectory forecasting model:

$$P_{\theta}(a_{t:t+T}^1, \dots, a_{t:t+T}^M | a_{t:t+T}^R, x_t)$$

future behavior of M agents Robot plan joint state:
 $x_t = [x_t^1, x_t^2, \dots, x_t^M, x_t^R]$



We also have a robot policy:

$$\Pi_{\phi}(x, P_{\theta}) = a_{t:t+T}^R$$

current state of world preds. ex. sequence of linear + angular ctrls

Deploy the robot for T timesteps in real world in N scenarios (ex. diff init conds, environments, cities, etc) and collect

$$\mathcal{D} = \left\{ \left(\underline{x_{0:T}}, \underline{a_{0:T}^{1:M}}, \underline{a_{0:T}^R} \right)_n \right\}_{n=1}^N$$

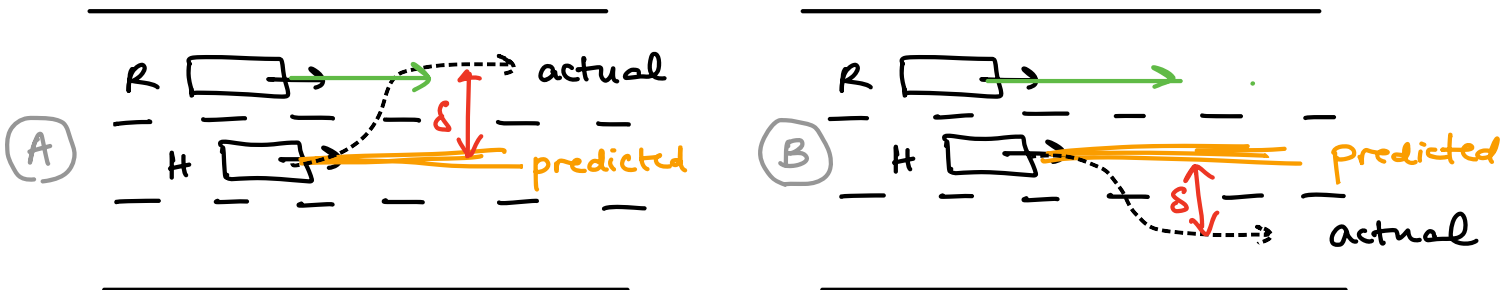
observed joint state traj. observed agent actions executed robot actions

This data contains real interactions btwn. the agents + robot.

We want to find $\mathcal{D}_{fail} \subset \mathcal{D}$ which have system-level

prediction failures: errors were poor predictions degraded perform.

Naive IDEA 1: Why not construct \mathcal{D}_{fail} as all scenarios where the predicted vs. real actions of $a_{0:T}^{1:M}$ deviated?



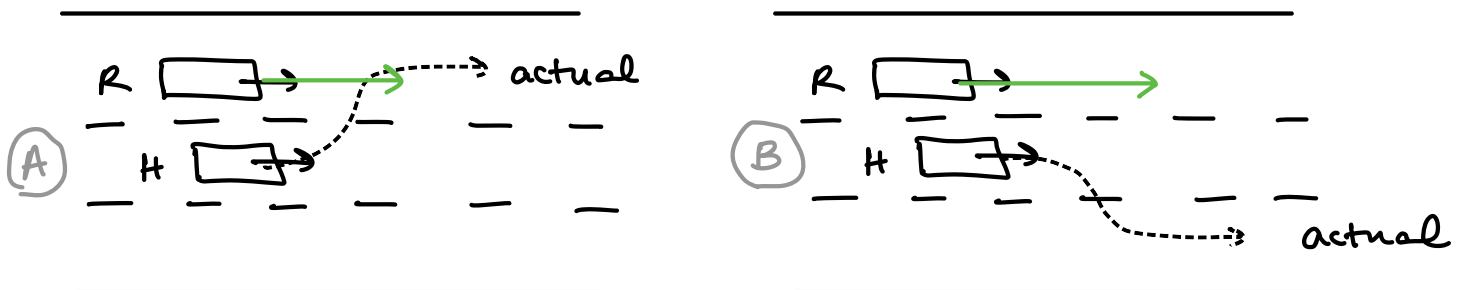
$$\delta = \| a_{0:T}^H - \hat{a}_{0:T}^H \|_2 \quad \equiv \quad \delta = \| a_{0:T}^H - \hat{a}_{0:T}^H \|_2$$

SAME L2 prediction error BUT (A) IS MUCH WORSE THAN (B)

KEY IDEA: system-level prediction errors are high regret:

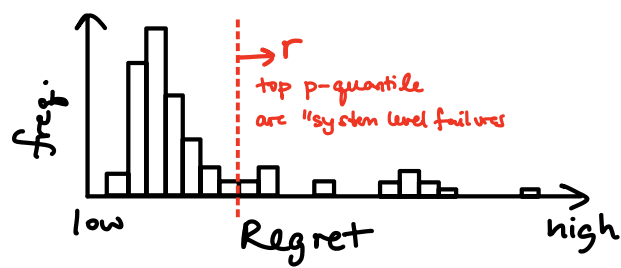
$$\text{Regret}_t := \max_{\bar{a}_{t:t+T}^R} R_\phi(\bar{a}_{t:t+T}^R, a_{t:t+T}^{1:M}, x_t) - R_\phi(a_{t:t+T}^R, a_{t:t+T}^{1:M}, x_t)$$

reward of the best robot decision in hindsight reward of executed robot decision AND observed agent behavior



$\max_{a^R} [R(a^R, \dots)] - R(\dots)$
 In hindsight, should have slowed down! \Rightarrow HIGH REGRET!

$\max_{a^R} [R(a^R, \dots)] - R(\dots)$
 In hindsight, would have still done \rightarrow \Rightarrow LOW REGRET!



$$\mathcal{D}_{fail} = \{ d = (a^{1:M}, a^R, x) : \text{Regret}(d) \text{ is in top } p\text{-quantile} \}$$

NOTE: high regret scenarios are more than collisions! ex. inefficiency is also high regret.

REPAIR

Identification is hard ... repair is even harder. But, there are two main techniques:

- ① re-training the faulty predictive model / component on the system-level failure data.
- ② monitor + fallback planning

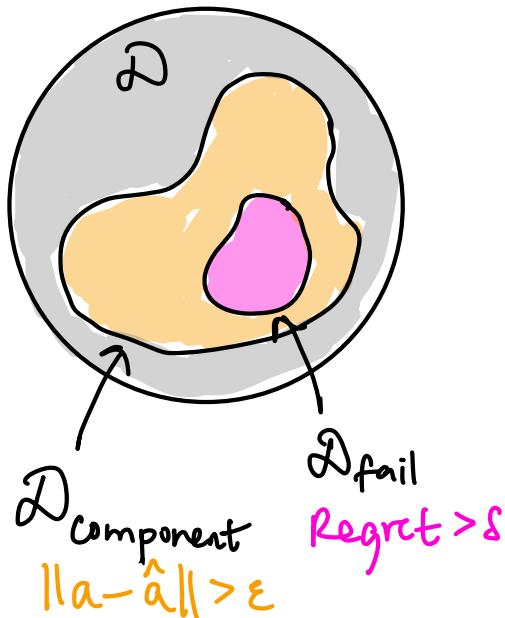
① In the perception case study, we can take all these:

$$BRT_I = \{ I : I = S(x; \theta), x \in BRT_x \}$$

and re-train the prediction model $I_t \rightarrow \square \rightarrow (\hat{p}_x, \hat{\theta})$

↳ in paper from above, re-training results in a 20% decrease in the BRT_I after re-training and 10% error reduction in the prediction error.

In the trajectory forecasting case-study, we can



If you fine-tune model on only D_{fail} (which is 77% less data than D !) reduces regret and collision cost during re-deployment by 53% & 65% respectively. BUT if you only fine-tune on LOW-REGRET DATA, you don't see statistically significant improvement.