

Last Time

- latent-space safety
- VLMs for monitoring & recovery

lecture 9

EACS S'25

Andrea Bajcsy

This Time:

- uncertainty quantification!

Announcement: midterm report due March 14<sup>th</sup> (Friday)

CREDIT: Notes inspired by Prof. Eric Nalisnick's lecture @ m<sup>2</sup>L

# Uncertainty Quantification for Predictive Models

So far, we have talked @ length about

safe decision-making / control

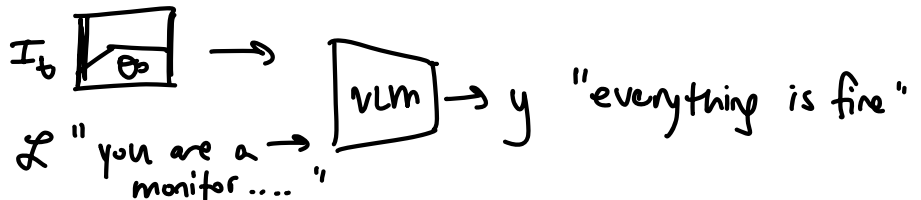
where safe  $\equiv$  constraint satisfaction

decision-making  $\equiv$  computing a policy / plan that makes sure that correct actions don't lead to future safety violations.

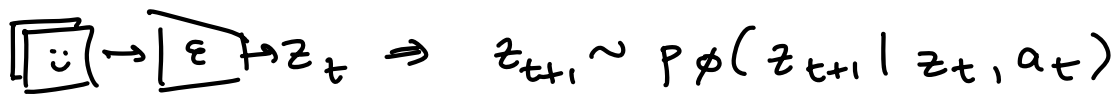
But, as we have pushed into the frontier, we have seen many more "components" or "models" that our decision-making depends on being "driven" by data.

These are all predictive models

ex. LLM / VLM which predicts next token



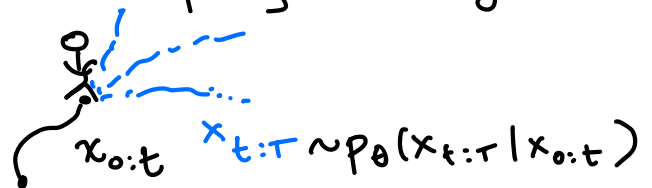
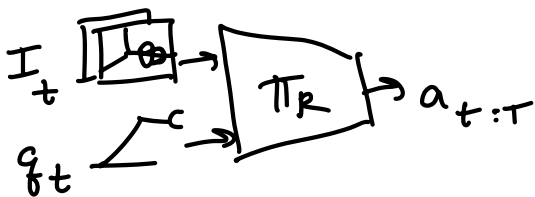
ex. world models which predict next latent state



ex. imitation policies

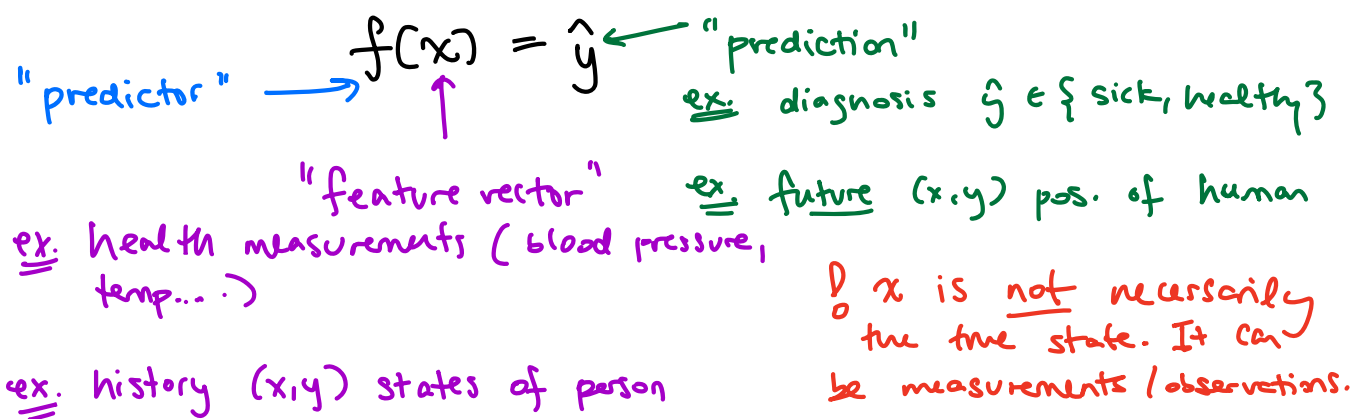
OR

human trajectory forecasting



Let's abstract these models, their architectures, etc. so we can unify our discussion.

In general, prediction problems look something like this.



Wait, what we want is to know how certain is  $f(\cdot)$  about its prediction! These predictive models will interact with "downstream" decision-making modules (e.g. doctor who gives you drugs, robot planner which looks @ predictions to take actions)

What we would like is something like a confidence statement

$$f(x) = \hat{y} \quad \underline{\text{AND}} \quad P(y = \hat{y} | x)$$

ex. 80% confident the person is sick.

Our goal is to know what our predictive models do AND do not know

## Two Types of Uncertainty

We first need to define what we mean by uncertainty.

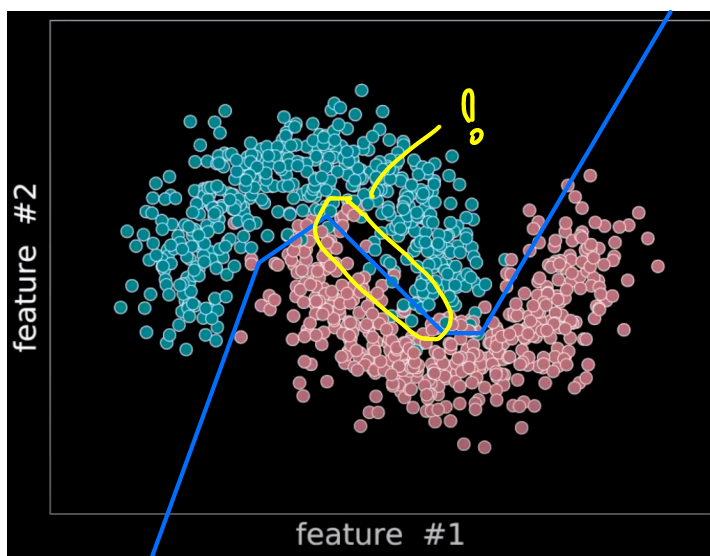
### 1) ALEATORIC

"lowest possible error rate of a classifier"

- fundamental, related to the Bayes error rate
- "irreducible" uncertainty even if you collect more data  
 $\Rightarrow$  only way around this is to collect more features

ex. train NN  
to learn classifier

$$f(x) = \hat{y}$$



← from Prof. Eric  
Nalisnick's lecture  
@ m2L.

high aleatoric uncertainty in yellow region b/c there is fundamental overlap btwn. the distributions (green in red ; red in green)

ex. suppose the true data is linear w/ Gaussian noise

$$y \sim \mathcal{N}(a + bx, \sigma^2)$$

The optimal estimator is linear regressor  $\hat{y} = \hat{a} + \hat{b}x$ . As we collect more data,  $\hat{a}, \hat{b} \rightarrow a, b$ . So the best error we could get is  $\sigma^2$ , the irreducible error (noise in the data itself).

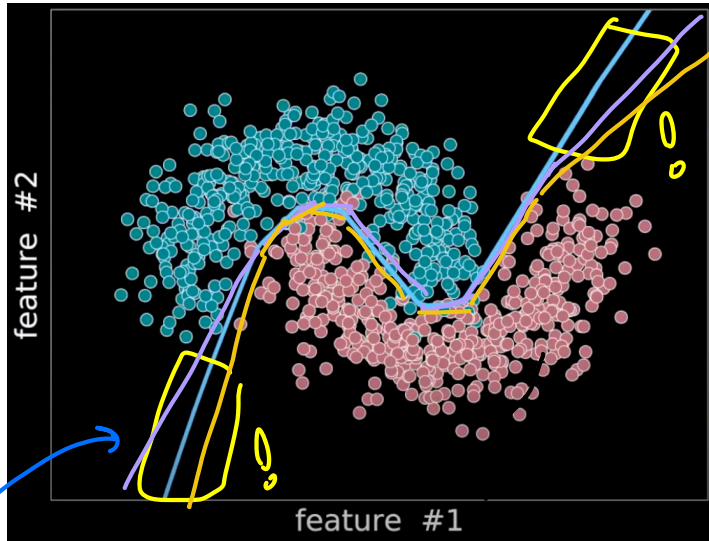
## 2. EPISTEMIC

- "easier" to deal with b/c it relates to a lack of data / experience
- always reduced by collecting more data!

! BUT in practice, hard!  
e.g. safety-critical data is less freq. and hard to get

ex. region of

high epistemic uncertainty is in yellow where our model makes predictions but has no data to be grounded with!



multiple training runs of NN with different rand. seeds.

It's important to understand these differences in uncertainty but in practice its very hard to know the diff!  
⇒ for most of these lectures, we will brush this distinction under the rug a bit; say uncertainty is high if either types are high.

### Notation & Assumptions

For these lectures, we will assume that there is a fixed, unknown distribution that generates data:

$$y \sim \mathbb{P}(y | x)$$

← features  
← true "labels"

We get to see a (finite) number of samples to form our training data:

$$\mathcal{D} := \{(x_i, y_i)\}_{i=1}^N$$

We fit a model to recover the ground-truth distribution:

$$f(x) := p(y|x) \approx \mathbb{P}(y|x)$$

## Modeling Paradigms

Broadly speaking, there are two modeling schools of thought:

frequentism and Bayesianism. An intuitive separation comes from where we model "randomness" as coming from.

1) FREQUENTISM: randomness comes from the data distribution.

What this translates to in terms of model learning is

learning the maximum likelihood estimator (MLE):

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \sum_{i=1}^N \log p(y_i | x_i; \theta)$$

$$\left. \begin{aligned} f_{\theta}(x) \\ := \\ p(y|x; \theta) \end{aligned} \right|$$

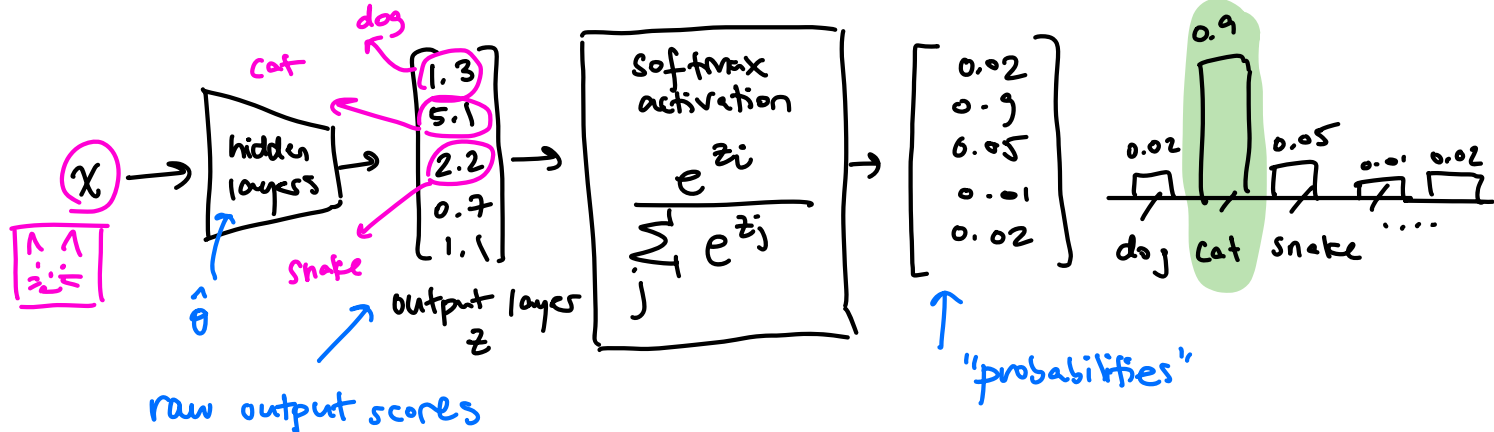
maximize the log likelihood of model parameters  $\theta$ .

FREQUENTIST IDEAL UQ — uncertainty quantific.

Ideally, under this paradigm, if I have a really big model and really big dataset, etc. we can quantify uncertainty by simply looking at the model probabilities

$$p(\hat{y}|x; \hat{\theta}) \approx \mathbb{P}(y = \hat{y} | x)$$

Here, UQ looks trivial! If I have classifier  $p(\hat{y}|x; \hat{\theta})$

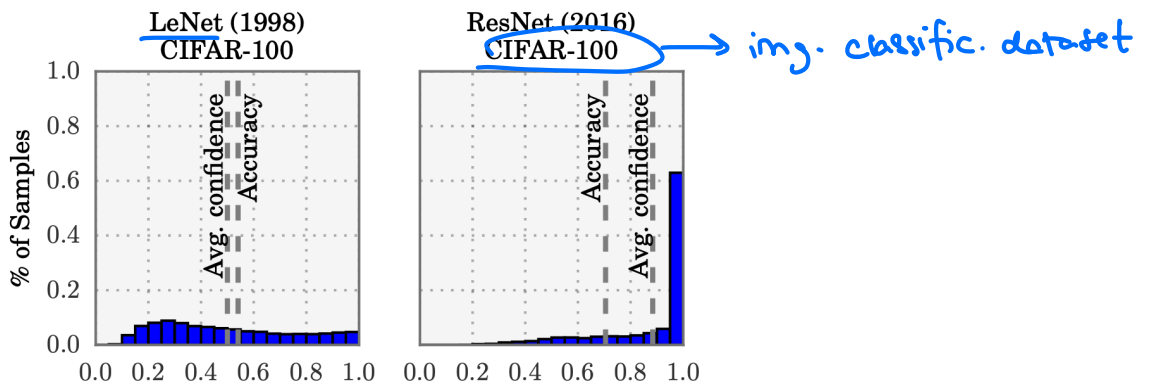


So are we "done"?

## Frequentist Learning: Limitations

In practice, as you may have experienced yourself, we can't usually rely on these probabilities directly.

ex. Guo et al. ICML 2017. "On Calib. of Modern NNs"



**LENET** is an old, smaller NN (5 layers). Showing softmax probabilities associated w/ each label

⇒ old model is "evenly distributed" confidence

⇒ avg. confidence matches accuracy (~50%)

**RESNET** bigger, more powerful model (110-layers)

⇒ new model has higher accuracy (~70%) BUT the average confidence is much HIGHER (~70%)

⚠ you can't read off softmax outputs b/c you can get overconfident estimate of how good you'd actually be.

2) BAYESIANISM: randomness is influenced by prior distribution over model parameters.

In Bayesian learning, we define a prior distribution  $p(\theta)$  over model parameters to "jump start" your learning - it can constrain your solutions to certain "plausible" solutions. You multiply your prior by your likelihood (this is where your model  $p(y_i | x_i; \theta)$  comes in) and then you normalize to get a posterior distribution that has been updated b/c you have seen data:

$$p(\theta | \mathcal{D}) = \frac{\overset{\text{prior}}{p(\theta)} \prod_{i=1}^N \overset{\text{likelihood}}{p(y_i | x_i; \theta)}}{\boxed{P(\mathcal{D})} := \int p(\tilde{\theta}) \prod_{i=1}^N p(y_i | x_i; \tilde{\theta}) d\tilde{\theta}}$$

"posterior" ↑

Normalizing constant is the hard part about Bayesian learn!  
when this is NN params, hard!

### IDEAL BAYESIAN QA

Assuming you could solve the normalizer, then given any new data point  $\tilde{x}$ , you can compute the posterior predictive distribution:

$$p(\tilde{y} | \tilde{x}, \mathcal{D}) = \int_{\theta} \overset{\text{pred. model}}{p(\tilde{y} | \tilde{x}; \theta)} \overset{\text{posterior}}{p(\theta | \mathcal{D})} d\theta$$

↑  
What you would use to make preds! B/c it has all the uncertainty of your model baked into it -- uncertainty over different models accounted for in posterior



Under (near) perfect learning, use post. pred. dist as your "ground-truth" probabilities

$$p(\tilde{y} | \tilde{x}, \mathcal{D}) \approx P(y = \tilde{y} | x)$$

You can report confidence just like before...

### Bayesian learning: Limitations

Mostly computational, integrating over params hard for NNs, and so computing normalizer or posterior pred. dist is hard.