



CSE598: Operationalizing Deep Learning: A Sociotechnical Perspective

Ransalu Senanayake

School of Computing and Augmented Intelligence
Arizona State University






WAYMO

R

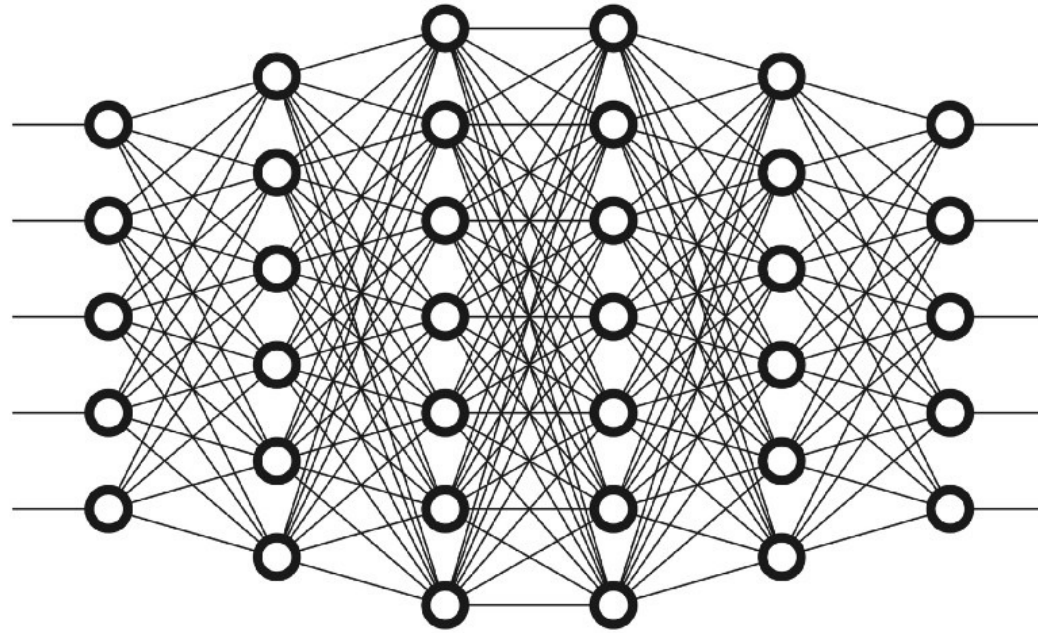
IVÁN BARRAGAN

AUTOPILOT DIDN'T REACT OR GIVE ANY WARNINGS, SO IVÁN TOOK OVER AND APPLIED THE BRAKES. HE DOESN'T THINK AUTOPILOT HAD A CHANCE TO REACT, AS THE ACCIDENT WAS TOO FAR IN THE DISTANCE.



KAPWING

We Don't Know



**Because We Don't Know How
Black Box ML Models Work**



**Chandler, Arizona: Waymo
Collision**



**Tempe, Arizona: Uber Hit And
Run**

**Mountain View, California:
Tesla Model X Crash**

Driver killed in self-driving car accident for first time

Nation Jun 30, 2016 5:24 PM EST

“The high ride height of the trailer of the truck combined with its positioning across the road and the rare circumstances of the impact caused the Model S to pass under the trailer” - Tesla



How Do we Train Deep Neural Networks?

Discriminative model vs. Generative model

$$p(y|x)$$

$$p(y, x)$$

Dataset

Training

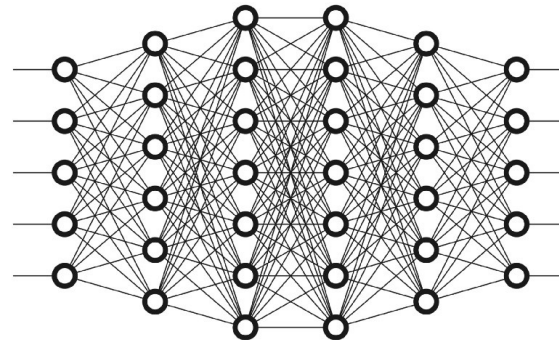
Testing



Predicted: Wolf
True: Wolf

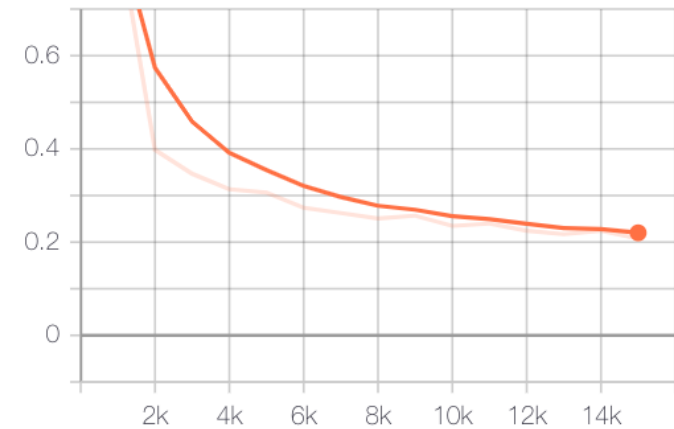


Predicted: Husky
True: Husky



Training with backpropagation
to minimize the error

training_loss





Predicted: **Wolf**
True: **Wolf**



Predicted: **Husky**
True: **Husky**



Predicted: **Husky**
True: **Husky**



Predicted: **Wolf**
True: **Wolf**



Predicted: **Wolf**
True: **Wolf**



Predicted: **Wolf**
True: **Wolf**



Predicted: **Husky**
True: **Wolf**



Predicted: **Wolf**
True: **Wolf**

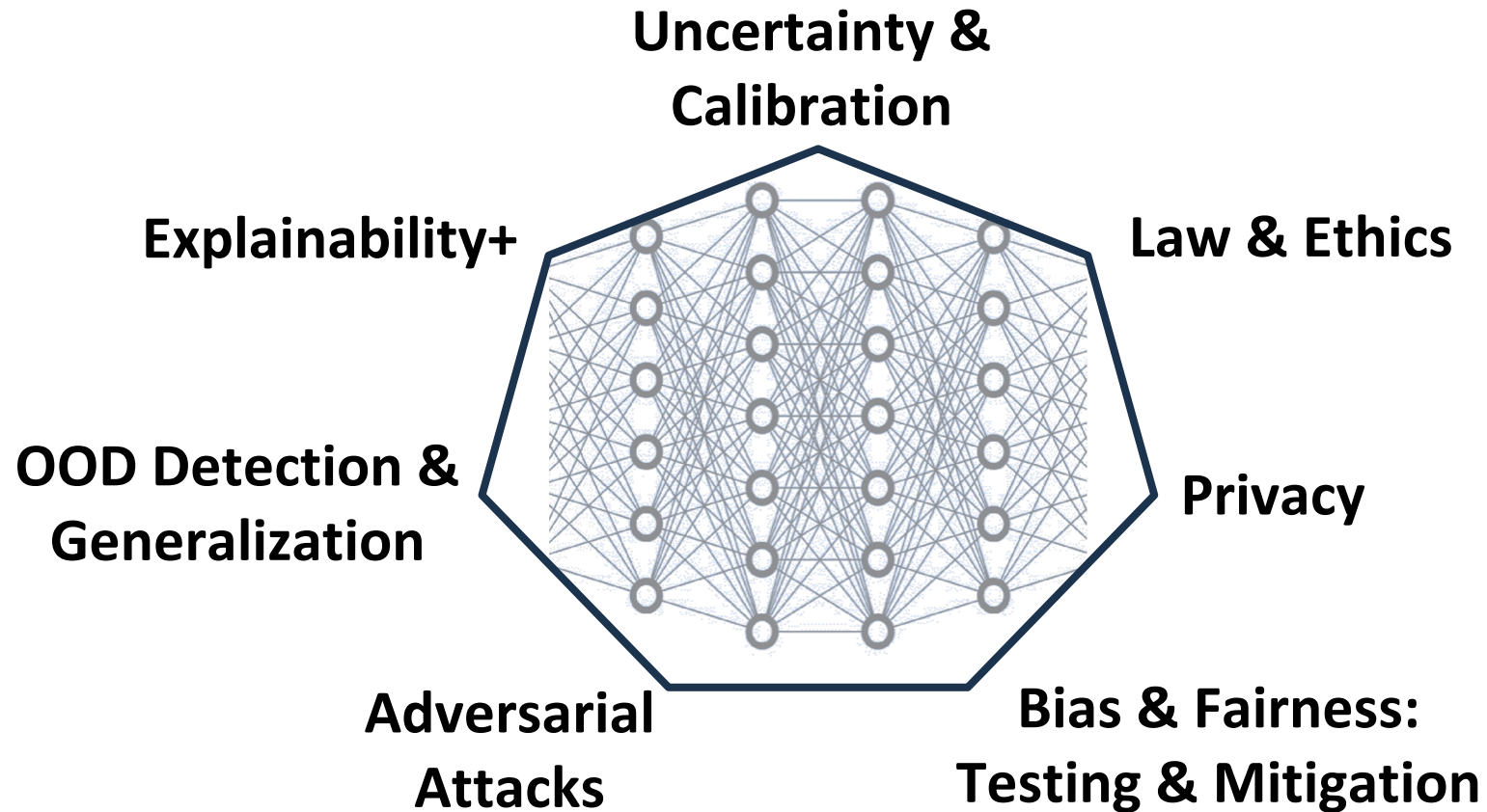


Predicted: **Wolf**
True: **Husky**



Predicted: **Husky**
True: **Husky**

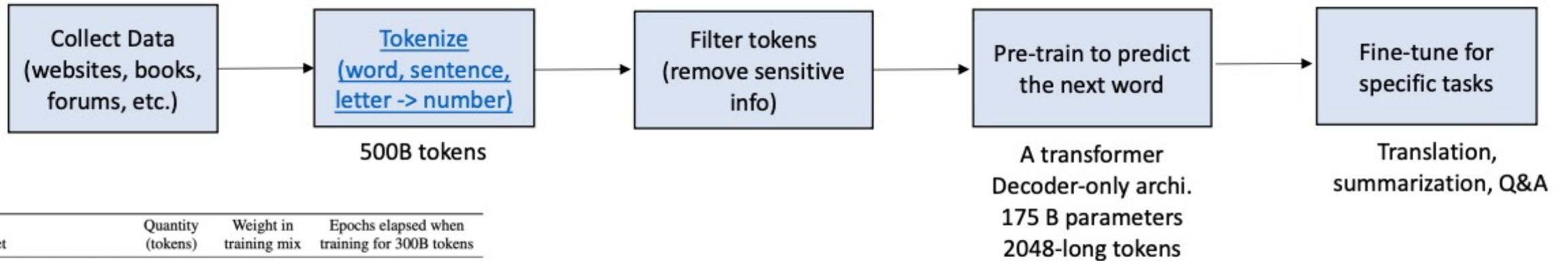
Is That Enough?



Application 1 (NLP)

Sentence Completion & Review Generation

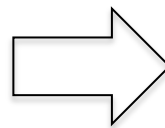
How Has GPT3 Been Trained?



Dataset	Quantity (tokens)	Weight in training mix	Epochs elapsed when training for 300B tokens
Common Crawl (filtered)	410 billion	60%	0.44
WebText2	19 billion	22%	2.9
Books1	12 billion	8%	1.9
Books2	55 billion	8%	0.43
Wikipedia	3 billion	3%	3.4



Social bias
(historical bias, life bias, etc.)



Algorithmic bias
(dataset bias, model bias, etc.)

Generative models in Vision

Application 2 (Robotics)

Perception to Decision-Making

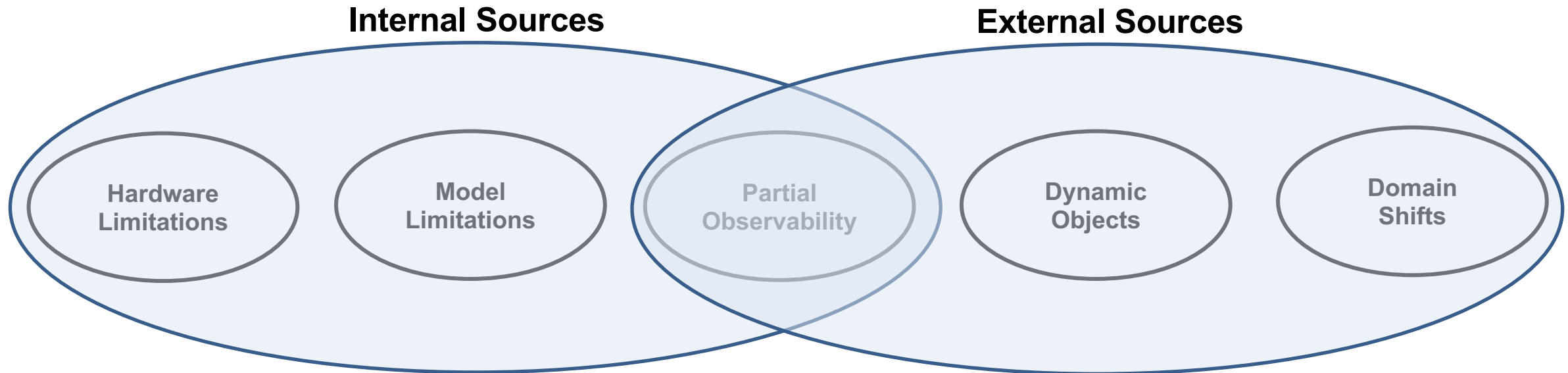


Сайт: vk.com/road

"Дорожные войны"

<http://vk.com/road>

Sources of Uncertainty



Sources of Uncertainty

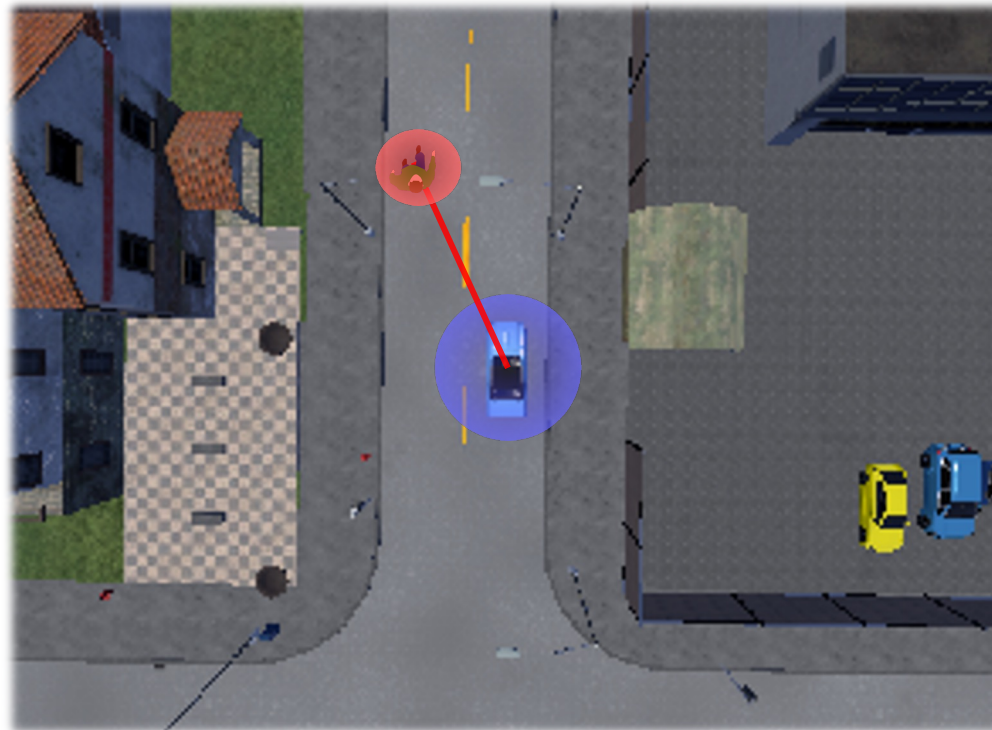
Hardware
Limitations

Model
Limitations

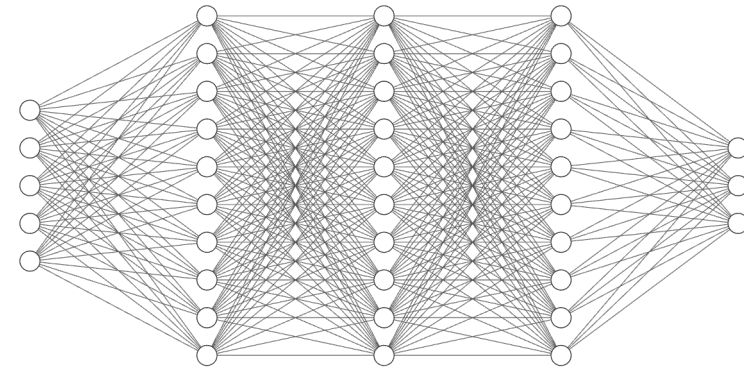
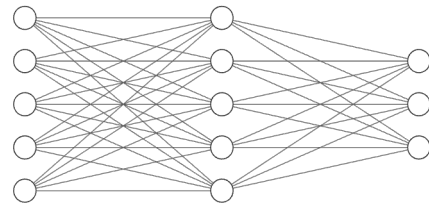
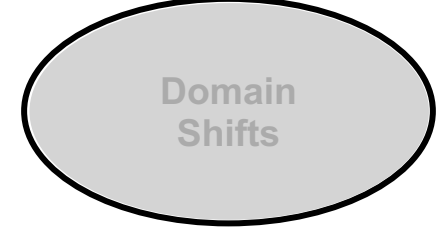
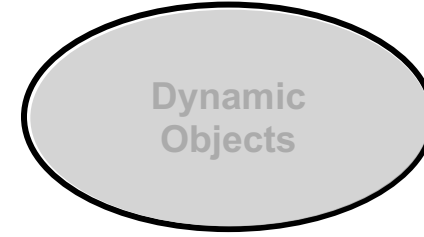
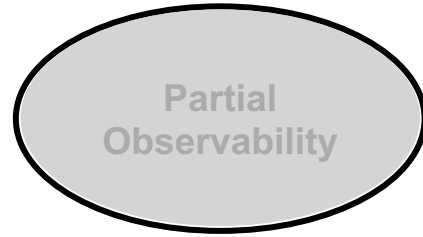
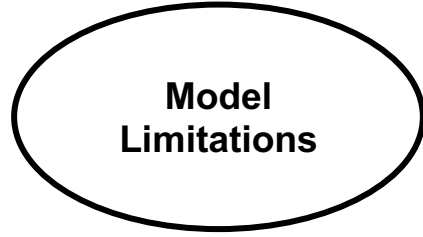
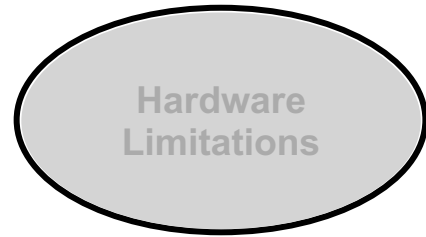
Partial
Observability

Dynamic
Objects

Domain
Shifts



Sources of Uncertainty



Sources of Uncertainty

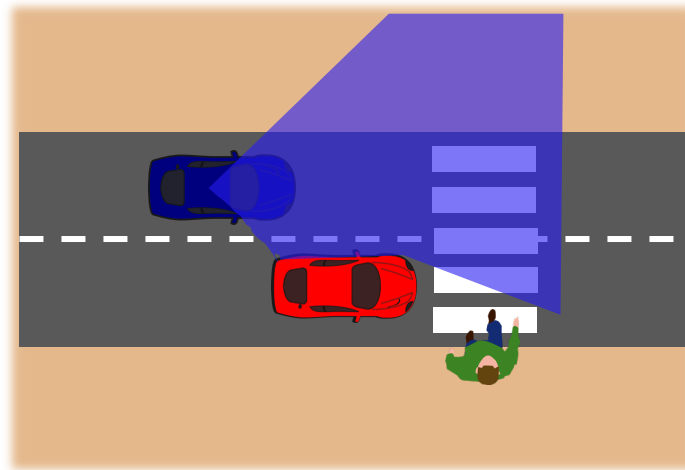
Hardware
Limitations

Model
Limitations

Partial
Observability

Dynamic
Objects

Domain
Shifts



Sources of Uncertainty

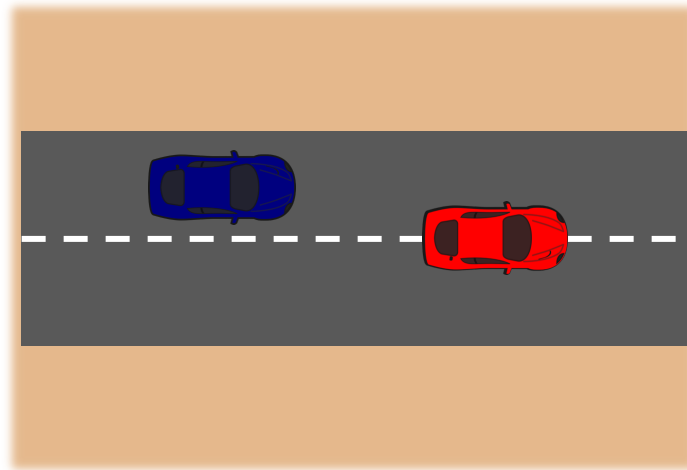
Hardware
Limitations

Model
Limitations

Partial
Observability

**Dynamic
Objects**

Domain
Shifts



Sources of Uncertainty

Hardware
Limitations

Model
Limitations

Partial
Observability

Dynamic
Objects

Domain
Shifts

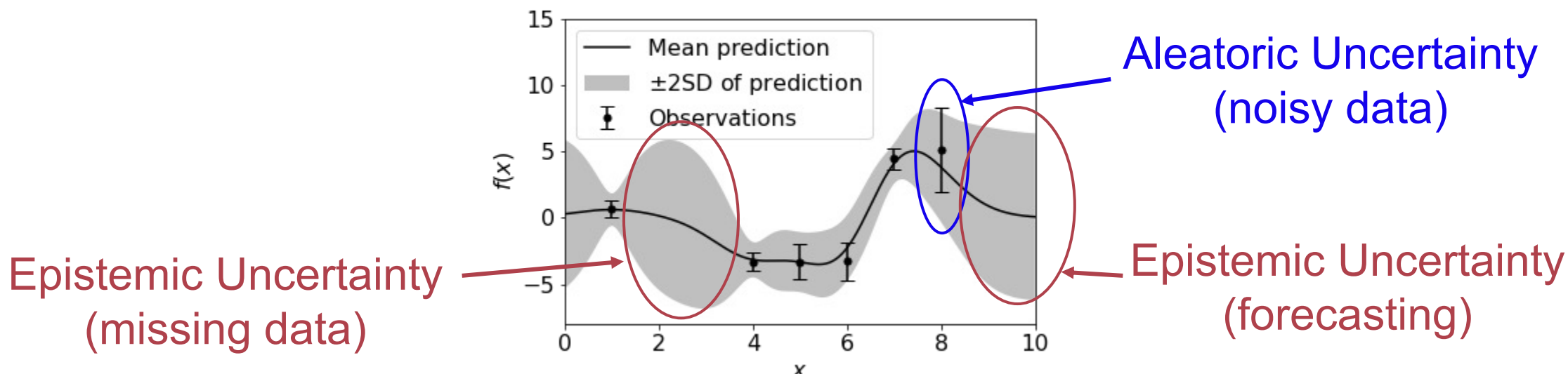
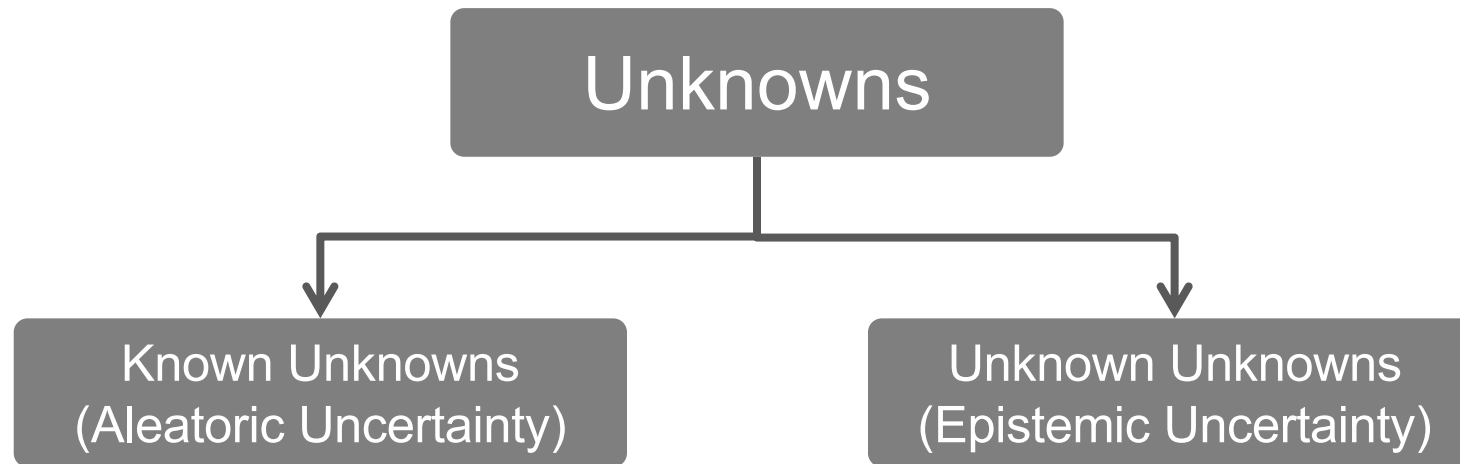


During training

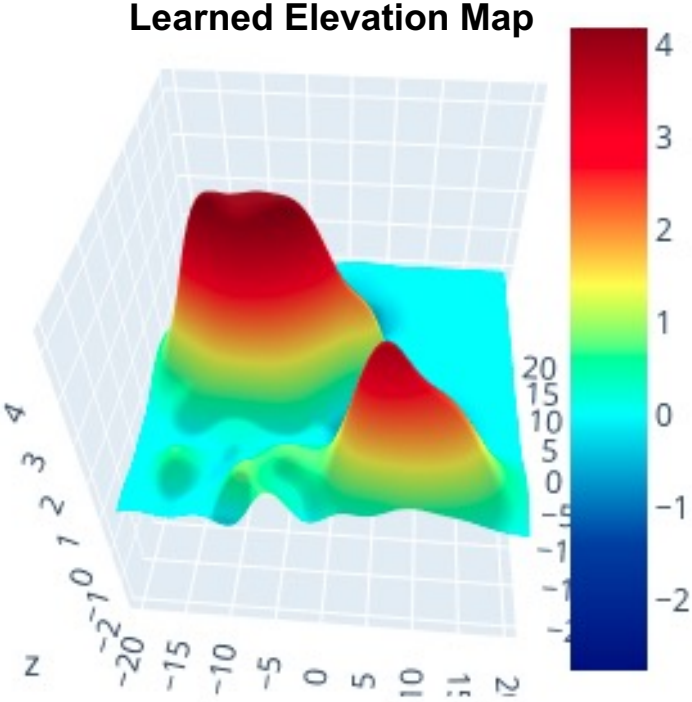
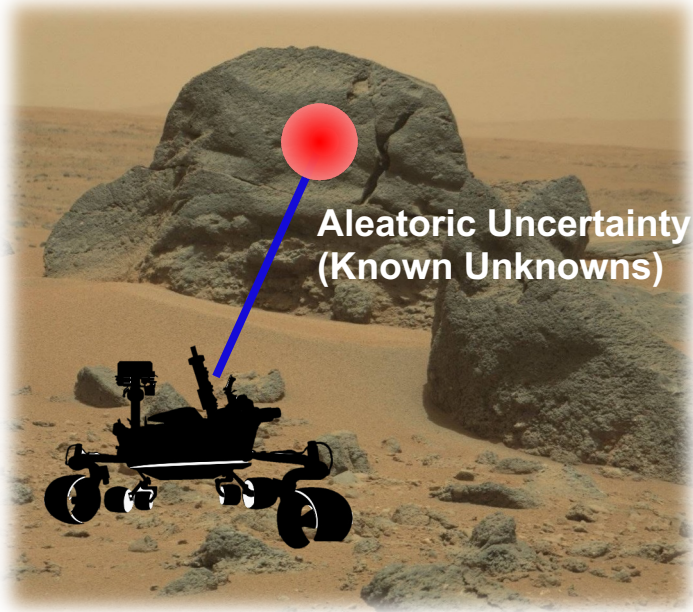


At deployment

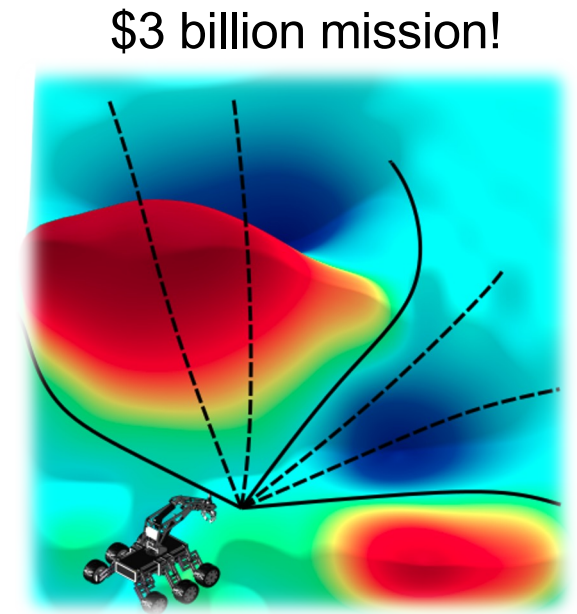
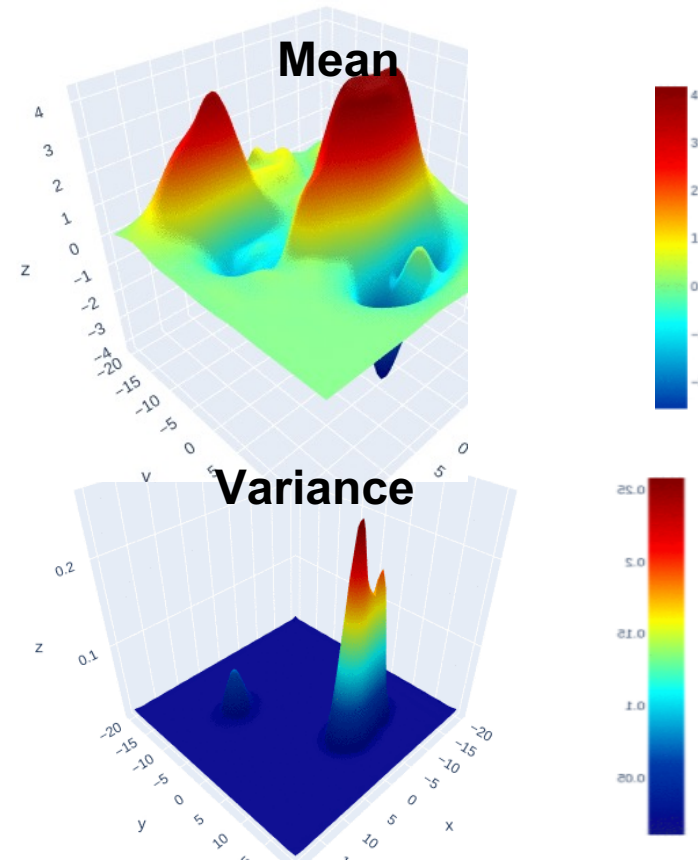
Types of Uncertainties



Known Unknowns (a.k.a. Aleatoric Uncertainty)

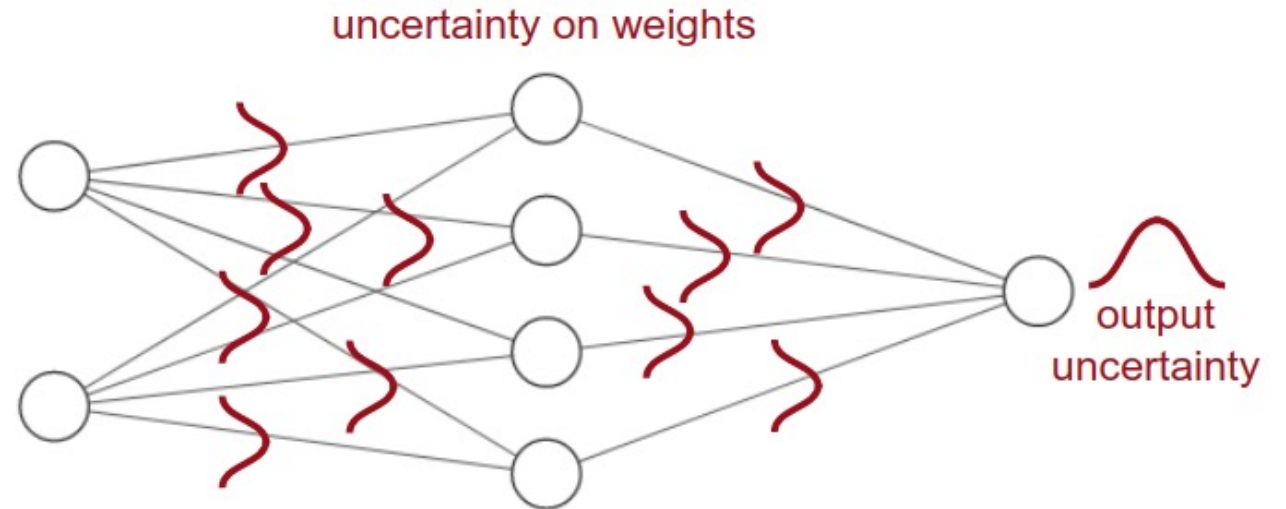


Unknown Unknowns (a.k.a. Epistemic Uncertainty)



Unknown Unknowns (a.k.a. Epistemic Uncertainty)

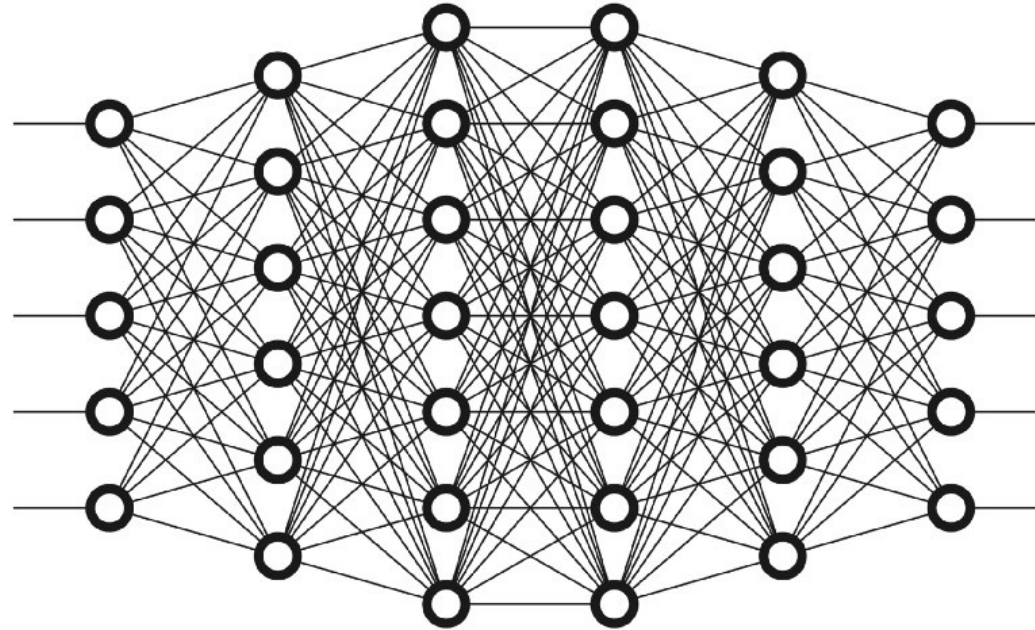
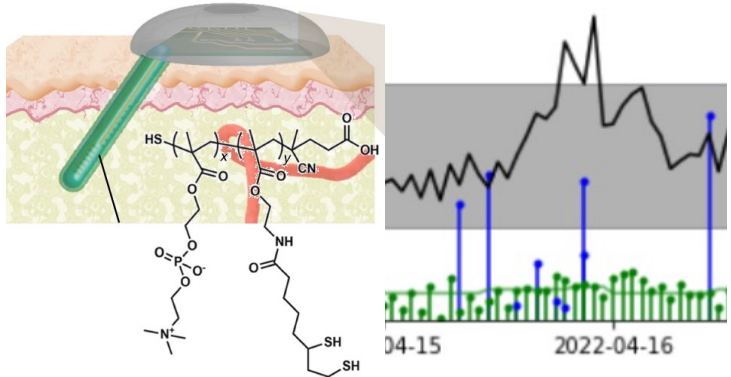
$$\hat{y} = \underbrace{\hat{w}_1 x + \hat{w}_0}_{\hat{f}_w(x)} + \epsilon$$



**It's okay to not know something.
However, we ought to know what we
don't know.**

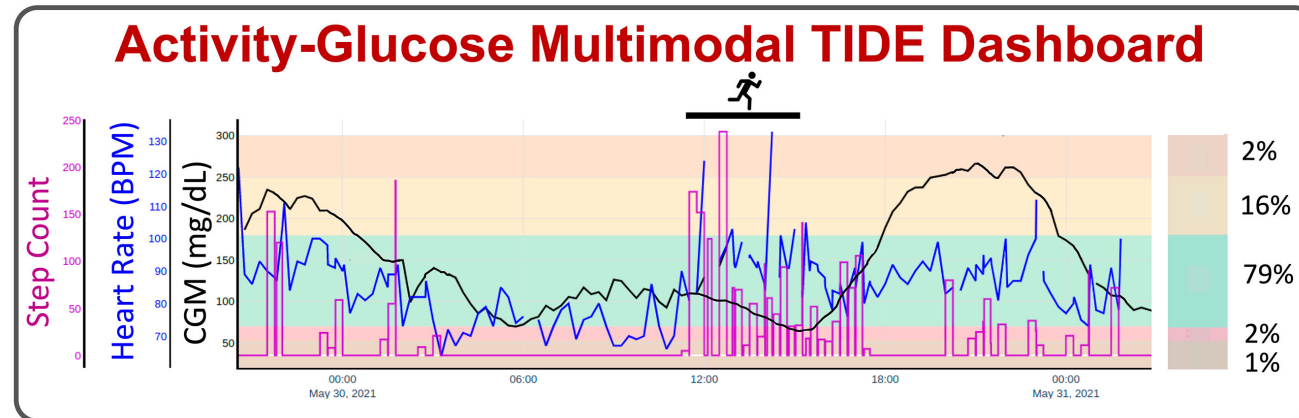
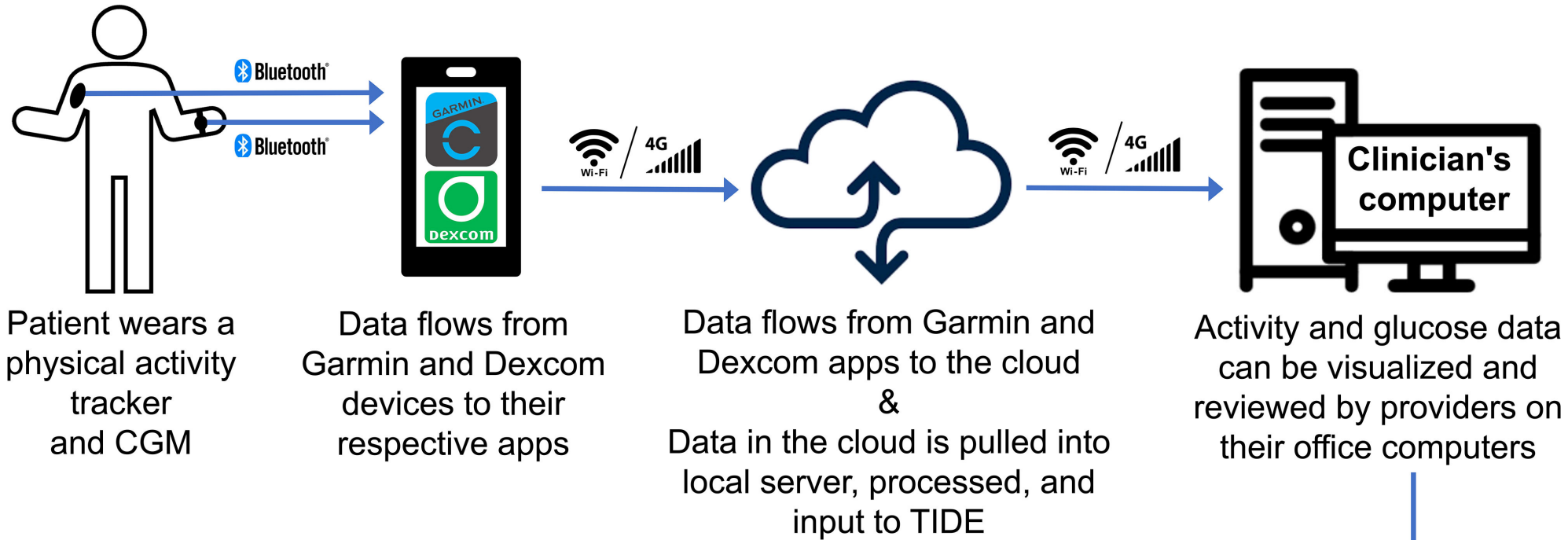
Application 3 (Healthcare)

Personalized Healthcare



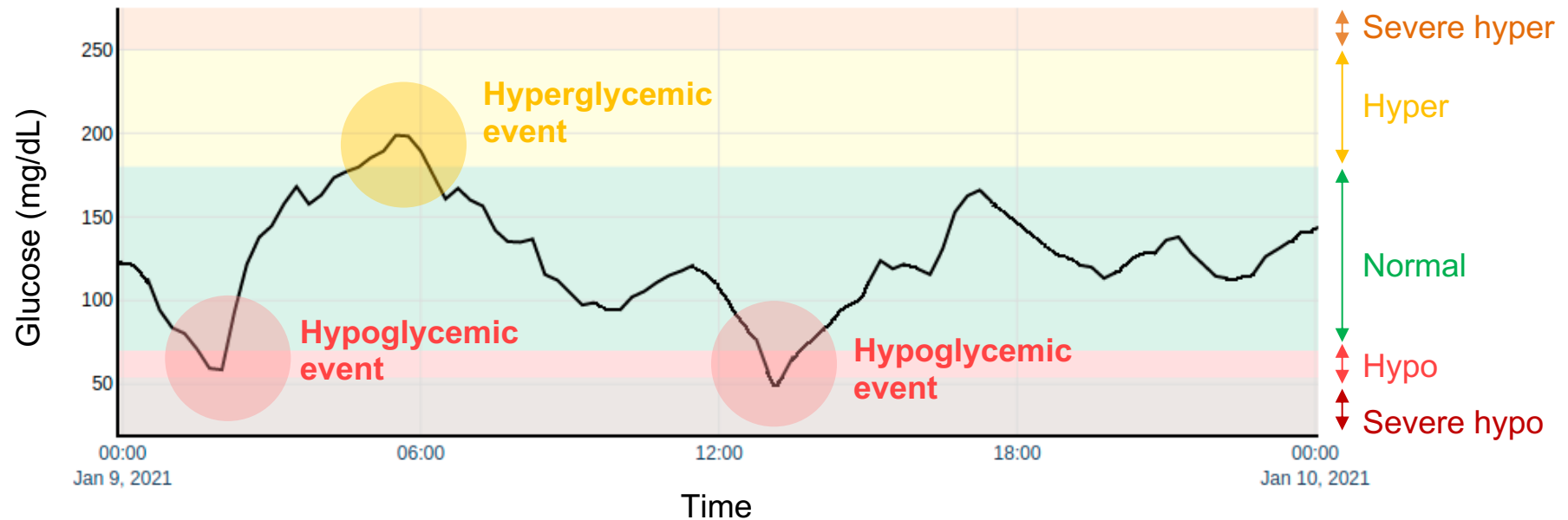
Output:
This patient has hypoglycemia

Clinicians need to understand!

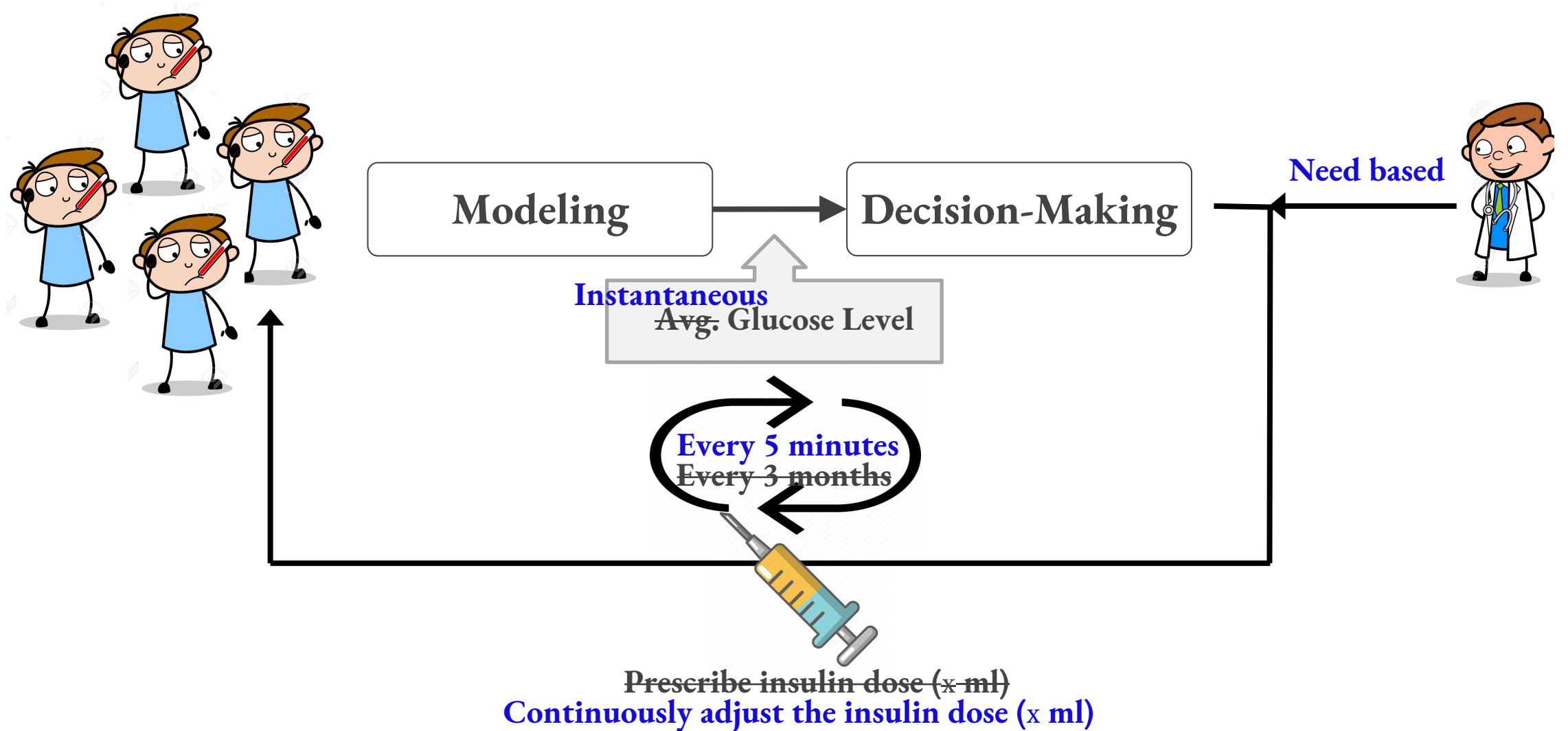


Type 1 Diabetes

The 3-month avg glucose (Hb1Ac) misses important events



Challenges



Training...

How Do we Train Deep Neural Networks?

Discriminative model vs. Generative model

$$p(y|x)$$

$$p(y, x)$$

Optimization

Training with backpropagation to minimize the error

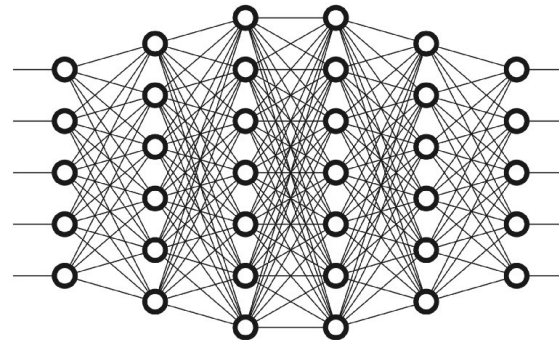
Dataset

Training

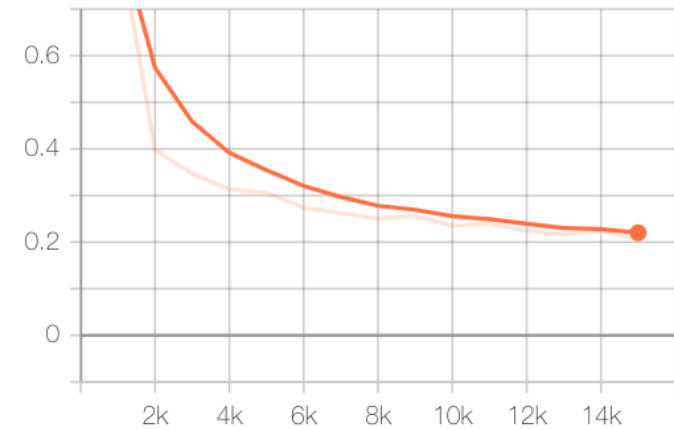
Testing

Data

Model



training_loss



Evaluation

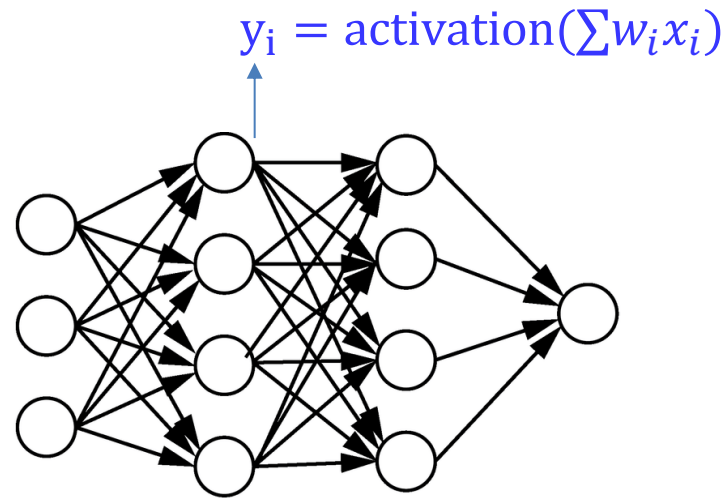


Predicted: Wolf
True: Wolf



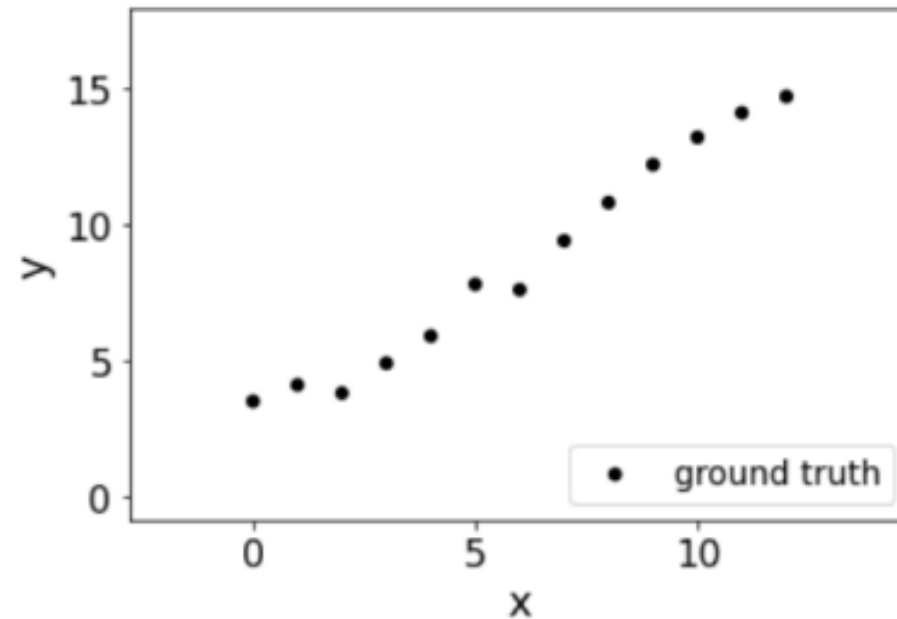
Predicted: Husky
True: Husky

Why Does a Neural Network Work So Well?



- 1) (Non)linear regression of (non)linear regressions view
- 2) Manifold view

Linear Regression: Problem



Objective: Predict outputs for unknown inputs x_q , given the training data $\{\mathbf{x}_i, y_i\}_{i=1}^N$.

$$y = f(\mathbf{x}) + e, \quad e \sim \mathcal{N}(0, \sigma^2)$$

Linear Regression: Model

$$y = \theta_0 + \theta_1 x^{(1)} + \theta_2 x^{(2)} + \theta_3 x^{(3)} + \dots + \theta_D x^{(D)} + e$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_1^{(D)} \\ 1 & x_2^{(1)} & \dots & x_2^{(D)} \\ 1 & x_3^{(1)} & \dots & x_3^{(D)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_N^{(1)} & \dots & x_N^{(D)} \end{bmatrix} \begin{bmatrix} \theta^{(0)} \\ \theta^{(1)} \\ \vdots \\ \theta^{(D)} \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ \vdots \\ e_N \end{bmatrix}$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \mathbf{e}, \quad \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

Linear Regression: Solution 1

Ordinary Least Square (OLS)

$$\theta_{OLS}^* = \underset{\theta}{\operatorname{argmin}} \|\mathbf{e}\|_2^2 = \underset{\theta}{\operatorname{argmin}} \|\mathbf{y} - X\theta\|_2^2$$

$$\begin{aligned} \|\mathbf{y} - X\theta\|_2^2 &= (\mathbf{y} - X\theta)^\top (\mathbf{y} - X\theta) \\ &= \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top X\theta - \theta^\top X^\top \mathbf{y} + \theta^\top X^\top X\theta \\ &= \mathbf{y}^\top \mathbf{y} - 2\theta^\top X^\top \mathbf{y} + \theta^\top X^\top X\theta \end{aligned}$$

$$(AB)^\top = B^\top A^\top$$

$$\theta^\top X^\top \mathbf{y} = \text{scalar} = (\theta^\top X^\top \mathbf{y})^\top = \mathbf{y}^\top X\theta$$

$$\frac{\partial}{\partial \theta} \|\mathbf{y} - X\theta\|_2^2 = -2X^\top \mathbf{y} + 2X^\top X\theta$$

By setting the derivative to zero,

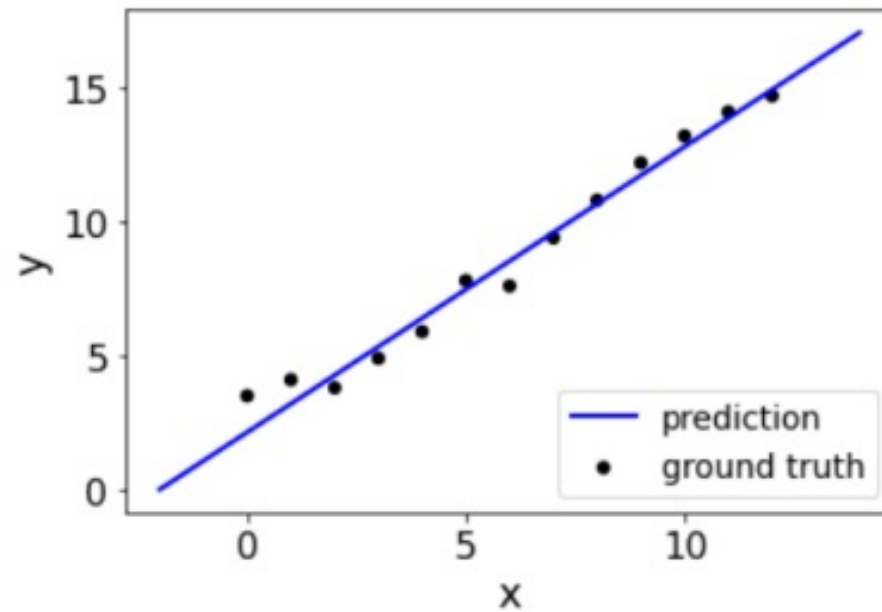
$$-2X^\top \mathbf{y} + 2X^\top X\theta_{OLS}^* = 0 \Rightarrow X^\top \mathbf{y} = X^\top X\theta_{OLS}^*$$

$$\boxed{\theta_{OLS}^* = (X^\top X)^{-1} X^\top \mathbf{y}}$$

Linear Regression: Predictions/Querying

Now, for unknown query input x_q , the output can be estimated as,

$$y_q = X_q \theta^*$$



$$y_q = \underbrace{\theta_0}_{=2.4} + \underbrace{\theta_1}_{=1.03} x_q$$

Linear Regression: Solution 2

Maximum Likelihood Estimate (MLE)

Alternatively, the same θ^* can be obtained by maximizing the likelihood,

$$\begin{aligned}\theta_{ML}^* &= \arg \max_{\theta} p(\mathbf{y}|\mathbf{X}, \theta) \quad (\text{maximize likelihood}) \\ &= \arg \max_{\theta} \log p(\mathbf{y}|\mathbf{X}, \theta) \\ &= \arg \min_{\theta} -\log p(\mathbf{y}|\mathbf{X}, \theta) \quad (\text{minimize negative log likelihood (NLL)}) \\ &= \arg \min_{\theta} (\mathbf{y} - \mathbf{X}\theta)^\top (\mathbf{y} - \mathbf{X}\theta) \\ &= \arg \min_{\theta} \|\mathbf{y} - \mathbf{X}\theta\|_2^2 \\ &= \theta_{OLS}^*\end{aligned}$$

$$p(\mathbf{y}|\mathbf{X}, \theta) = \frac{1}{\sqrt{(2\pi)^{D+1} |\sigma^2 I|}} \exp \left(-\frac{1}{2} (\mathbf{y} - \mathbf{X}\theta)^\top (\sigma^2 I)^{-1} (\mathbf{y} - \mathbf{X}\theta) \right)$$

Linear Regression: Solution 3

Maximum-A-Posterior (MAP)

Alternatively, θ^* can be obtained by maximizing the the posterior distribution (i.e. the mode of the posterior distribution),

$$\begin{aligned}\theta_{MAP}^* &= \arg \max_{\theta} p(\theta|X, \mathbf{y}) \quad \textit{maximum-a-posteriori (MAP)} \\ &= \arg \max_{\theta} \frac{p(\mathbf{y}|X, \theta)p(\theta)}{p(\mathbf{y})} \\ &= \arg \max_{\theta} p(\mathbf{y}|X, \theta)p(\theta) \\ &= \arg \max_{\theta} \underbrace{\log p(\mathbf{y}|X, \theta)}_{\text{log likelihood}} + \underbrace{\log p(\theta)}_{\text{log prior}}\end{aligned}$$

$$\underbrace{p(\theta|X, \mathbf{y})}_{\text{posterior}} = \frac{\overbrace{p(\mathbf{y}|X, \theta)}^{\text{likelihood}} \times \overbrace{p(\theta)}^{\text{prior}}}{\underbrace{p(\mathbf{y})}_{\text{marginal likelihood}}}$$

Linear Regression: Solution 3

Maximum-A-Posterior (MAP)

- The log-prior acts as a regularizer (penalizer) and prevents over-fitting/multicollinearity
- If $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \sigma_0^2 I)$,

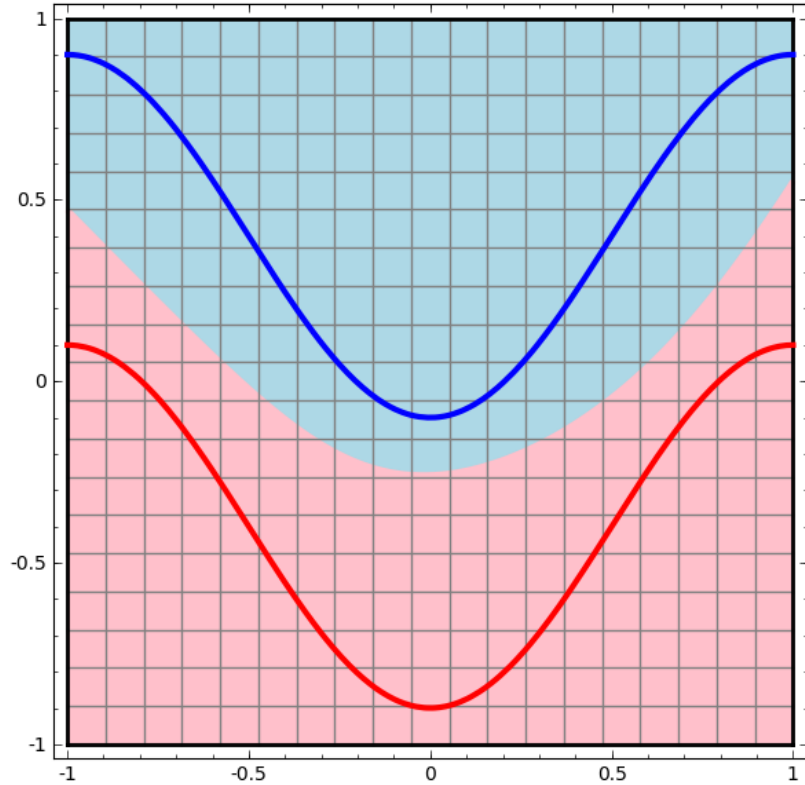
$$\boldsymbol{\theta}_{MAP}^* = (X^\top X + \lambda I)^{-1} X^\top \mathbf{y}, \quad \lambda = \sigma^2 / \sigma_0^2$$

- This is equivalent to ridge regression (Tikhonov or L2 regularization)
- λI term improves the numerical stability of the inversion

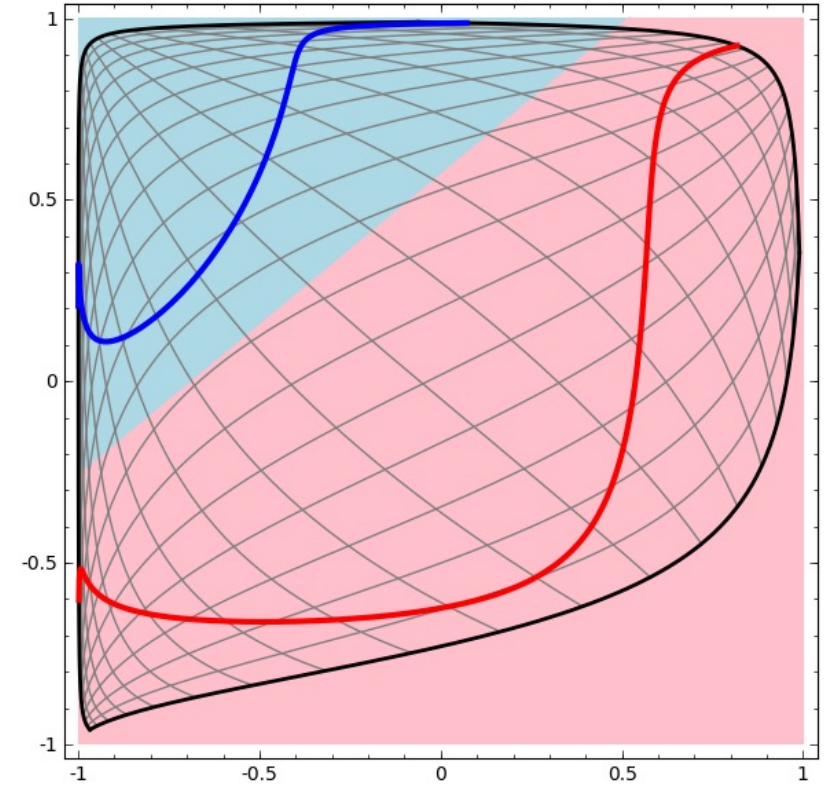
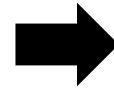
$$\boldsymbol{\theta}_{ML}^* = (X^\top X)^{-1} X^\top \mathbf{y}$$

Neural Networks: Manifolds View

<https://colah.github.io/posts/2014-03-NN-Manifolds-Topology/>



activation($\sum w_i x_i$)



Metrics

Confusion Matrix

		Predicted condition	
		Predicted Positive (PP)	Predicted Negative (PN)
Actual condition	Total population = P + N		
	Positive (P)	True positive (TP), hit	False negative (FN), type II error, miss, underestimation
	Negative (N)	False positive (FP), type I error, false alarm, overestimation	True negative (TN), correct rejection

Confusion Matrix

		Predicted condition			
		Predicted Positive (PP)	Predicted Negative (PN)		
Total population = P + N				Informedness, bookmaker informedness (BM) = TPR + TNR - 1	Prevalence threshold (PT) $= \frac{\sqrt{\text{TPR} \times \text{FPR}} - \text{FPR}}{\text{TPR} - \text{FPR}}$
Actual condition	Positive (P)	True positive (TP), hit	False negative (FN), type II error, miss, underestimation	True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power $= \frac{\text{TP}}{\text{P}} = 1 - \text{FNR}$	False negative rate (FNR), miss rate $= \frac{\text{FN}}{\text{P}} = 1 - \text{TPR}$
	Negative (N)	False positive (FP), type I error, false alarm, overestimation	True negative (TN), correct rejection	False positive rate (FPR), probability of false alarm, fall-out $= \frac{\text{FP}}{\text{N}} = 1 - \text{TNR}$	True negative rate (TNR), specificity (SPC), selectivity $= \frac{\text{TN}}{\text{N}} = 1 - \text{FPR}$
Prevalence $= \frac{\text{P}}{\text{P} + \text{N}}$		Positive predictive value (PPV), precision $= \frac{\text{TP}}{\text{PP}} = 1 - \text{FDR}$	False omission rate (FOR) $= \frac{\text{FN}}{\text{PN}} = 1 - \text{NPV}$	Positive likelihood ratio (LR+) $= \frac{\text{TPR}}{\text{FPR}}$	Negative likelihood ratio (LR-) $= \frac{\text{FNR}}{\text{TNR}}$
Accuracy (ACC) $= \frac{\text{TP} + \text{TN}}{\text{P} + \text{N}}$		False discovery rate (FDR) $= \frac{\text{FP}}{\text{PP}} = 1 - \text{PPV}$	Negative predictive value (NPV) = $\frac{\text{TN}}{\text{PN}}$ = 1 - FOR	Markedness (MK), deltaP (Δp) = PPV + NPV - 1	Diagnostic odds ratio (DOR) $= \frac{\text{LR}+}{\text{LR}-}$
Balanced accuracy (BA) $= \frac{\text{TPR} + \text{TNR}}{2}$		F_1 score $= \frac{2\text{PPV} \times \text{TPR}}{\text{PPV} + \text{TPR}} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$	Fowlkes–Mallows index (FM) $= \sqrt{\text{PPV} \times \text{TPR}}$	Matthews correlation coefficient (MCC) $= \frac{\sqrt{\text{TPR} \times \text{TNR} \times \text{PPV} \times \text{NPV}}}{-\sqrt{\text{FNR} \times \text{FPR} \times \text{FOR} \times \text{FDR}}}$	Threat score (TS), critical success index (CSI), Jaccard index $= \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}}$

Wikipedia article has the best summary table I've seen so far

Precision vs. Recall

- Precision (TP/TP+FP): Accuracy of the positive predictions.
- Recall/sensitivity (TP/TP+FN): Ability to capture/retrieve all positive samples. Useful when missing positives is a bad thing.
- F1 score: $2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$
- Sometimes working with precision or recall separately makes more sense than using F1 (e.g., in a medical diagnosis, we might need to focus more on maximizing precision)
- When there is an imbalance in the number of ground truth positive and negative samples, F1 score will be biased
- For different classification thresholds, we can plot the precision-recall curve. Generally, if we decrease the classification threshold (say, from 0.5 to 0.3), TP+FP go up (something remotely looks like a positive will now be a positive, whether correct or not). This will increase fall positives (because false positives typically lie around 0.5) which will then decrease precision.

Spam vs. ham
Cancer diagnosis

