

Measuring Performance

CS 474

Dr. Porter Jenkins

D&C 130: 18-19

18 Whatever principle of ^aintelligence we attain unto in this life, it will rise with us in the ^bresurrection.

19 And if a person gains more ^aknowledge and intelligence in this life through his ^bdiligence and obedience than another, he will have so much the ^cadvantage in the world to come.

Overfitting

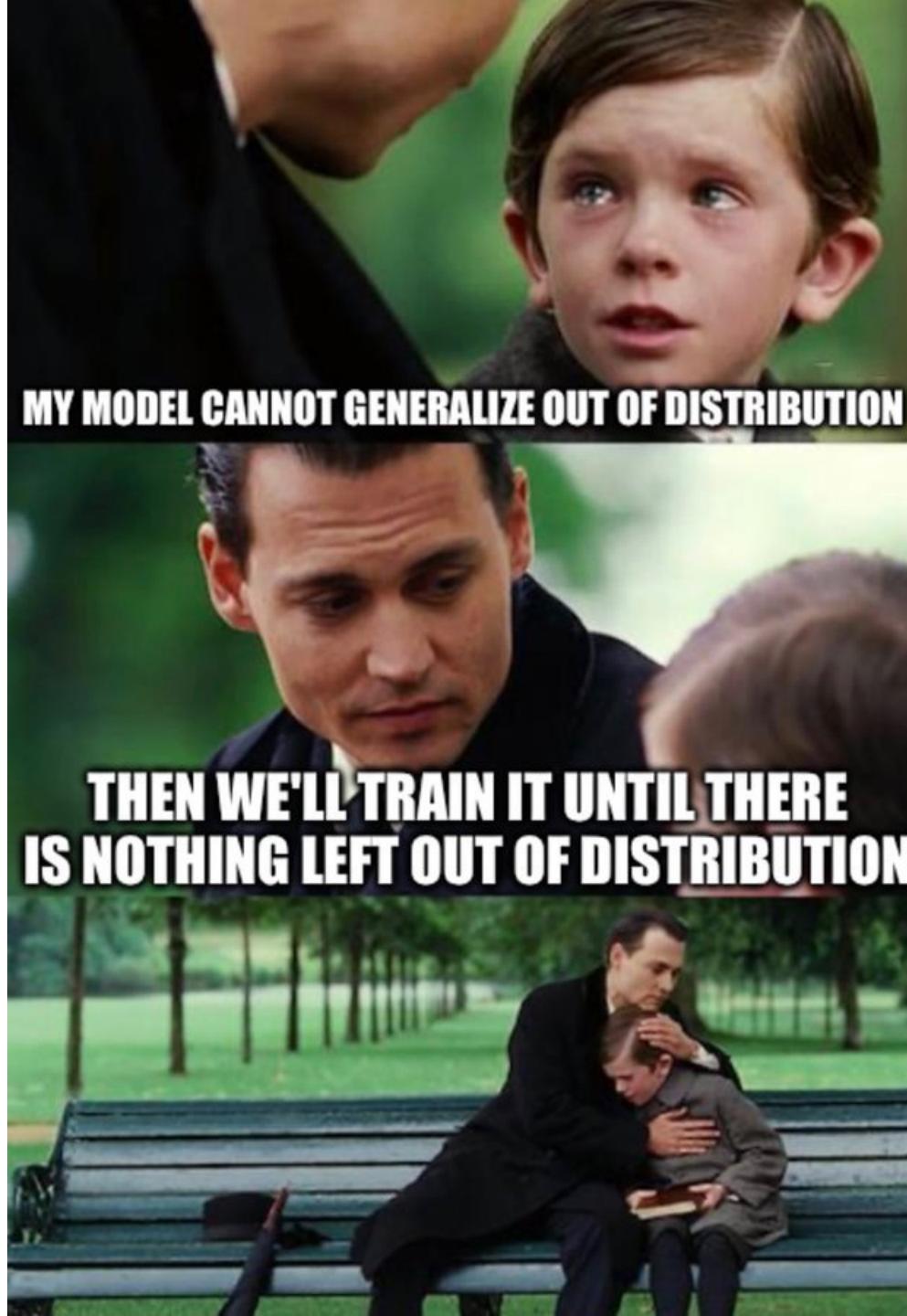
Overfitting

Phenomenon of machine learning algorithms to learn patterns specific to training data that do not generalize to unseen (test) data

In general, as we increase the capacity of the model, and broaden the family of possible functions we can represent, we increase the risk of overfitting

Never underestimate the ability of your neural network to overfit!!

Overfitting



Overfitting

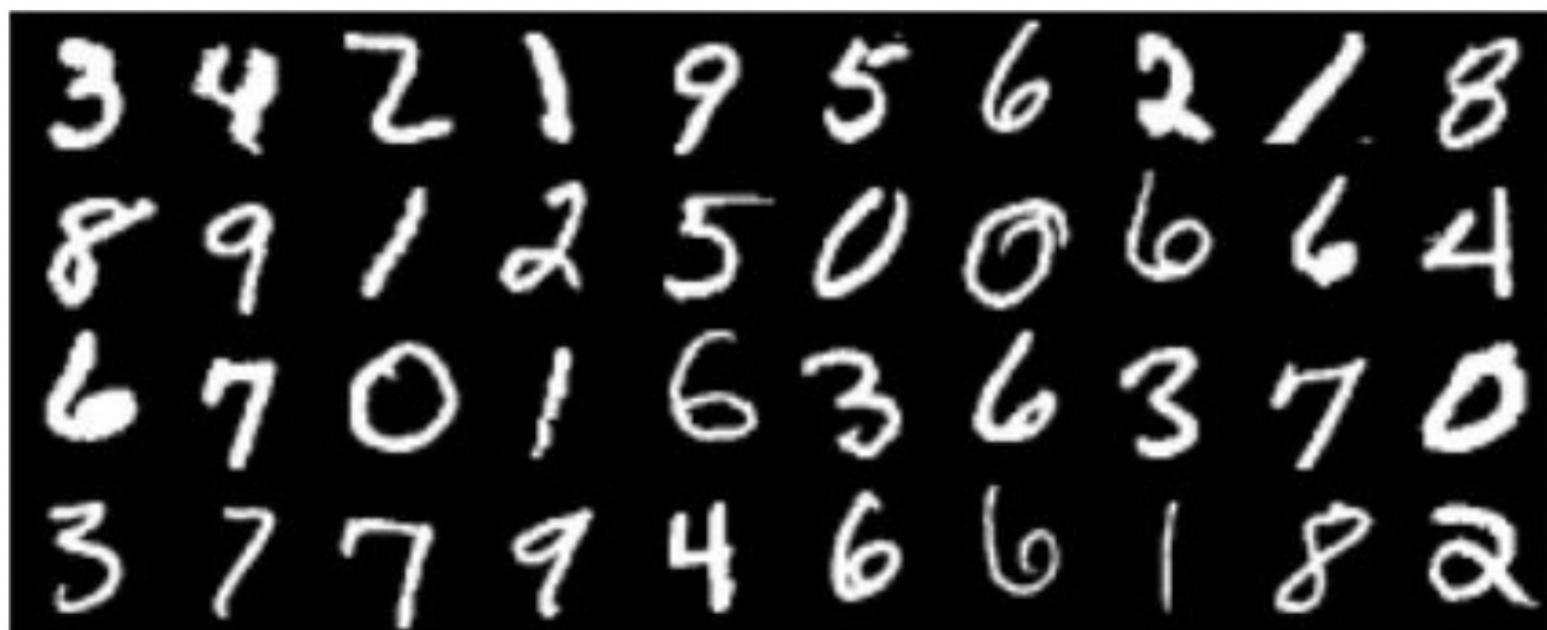
Demo:

https://colab.research.google.com/drive/1hBcLgVu8JBJR3_E-vOSoy63XBQTpeGO?authuser=2

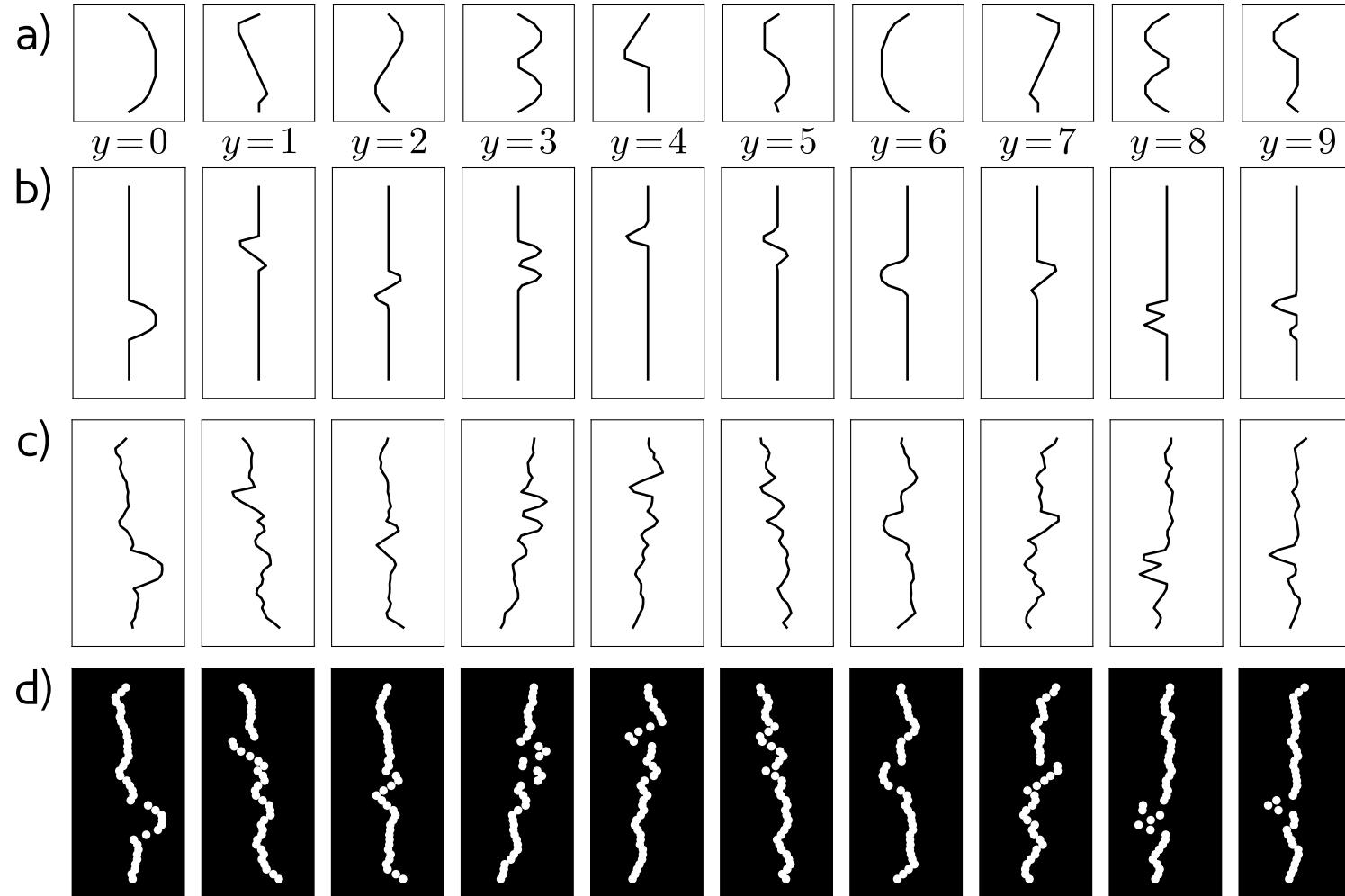
Measuring performance

- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Curse of dimensionality & weird properties of high dimensional space
- Choosing hyperparameters

MNIST Dataset



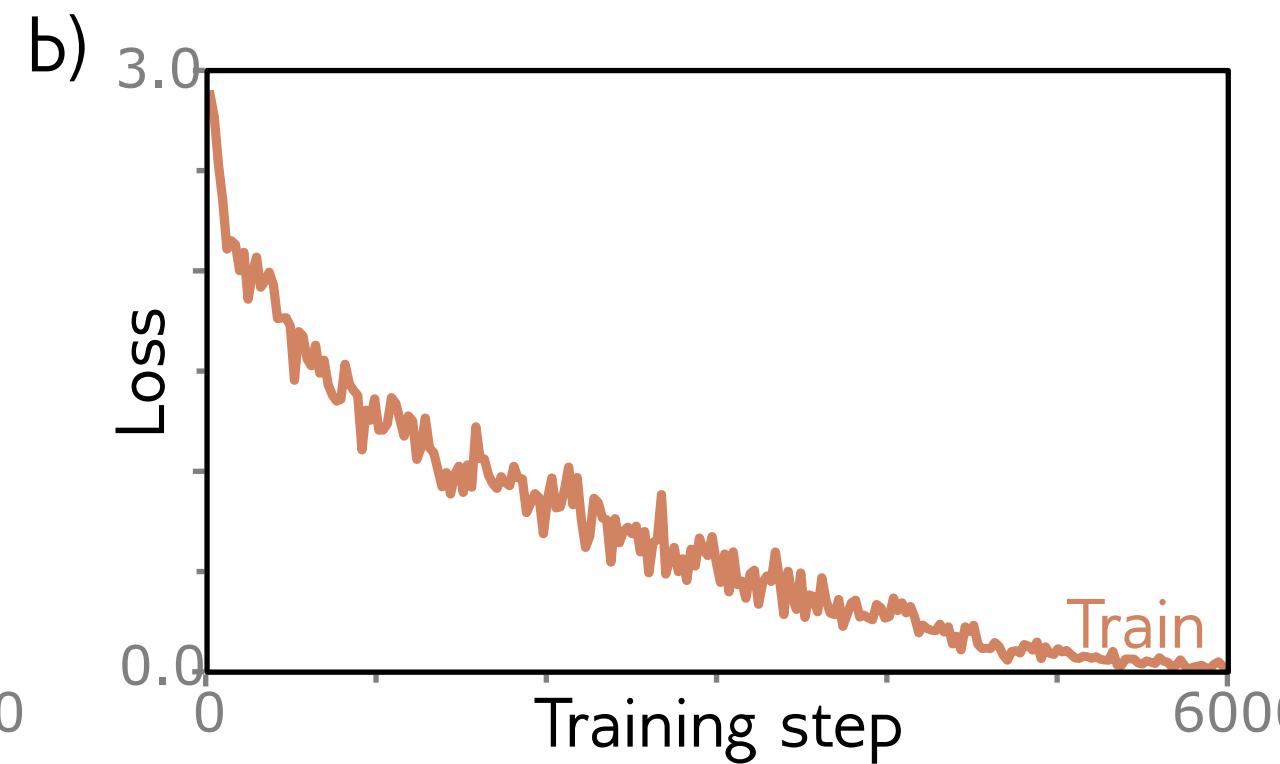
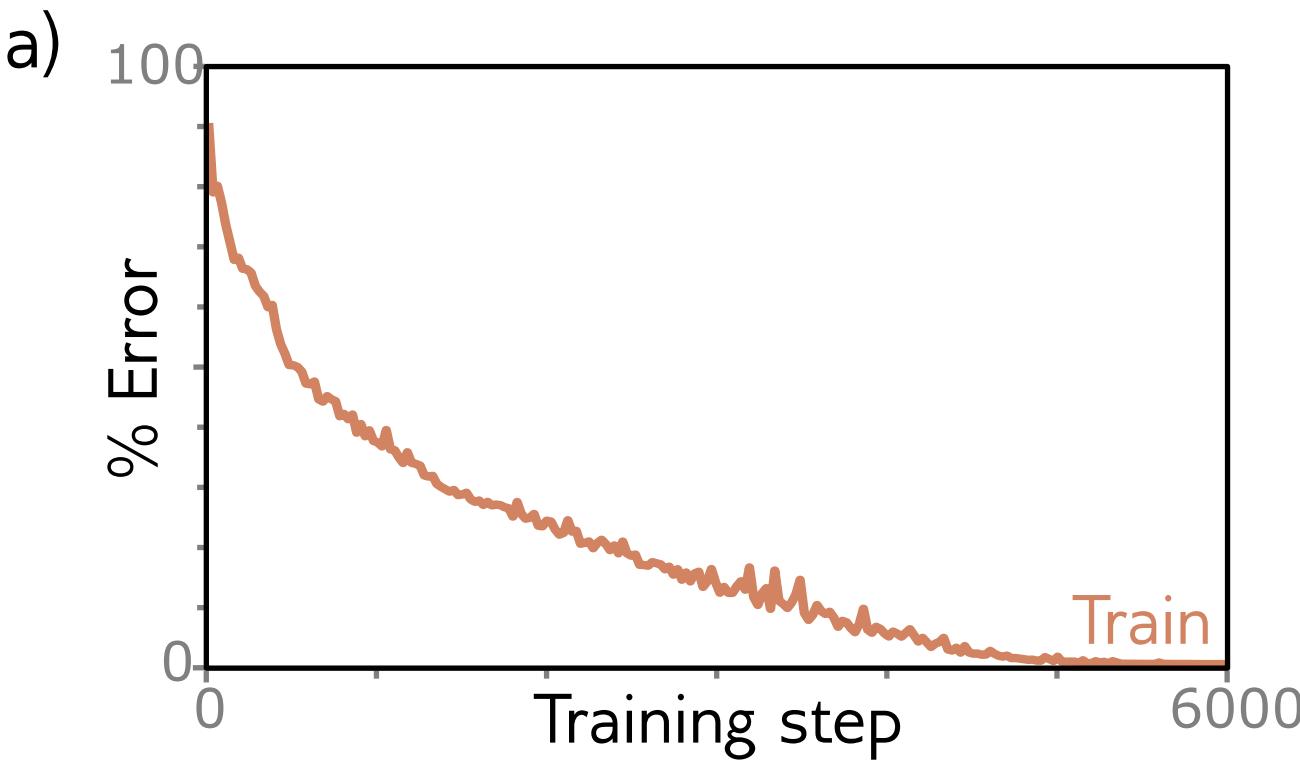
MNIST 1D Dataset



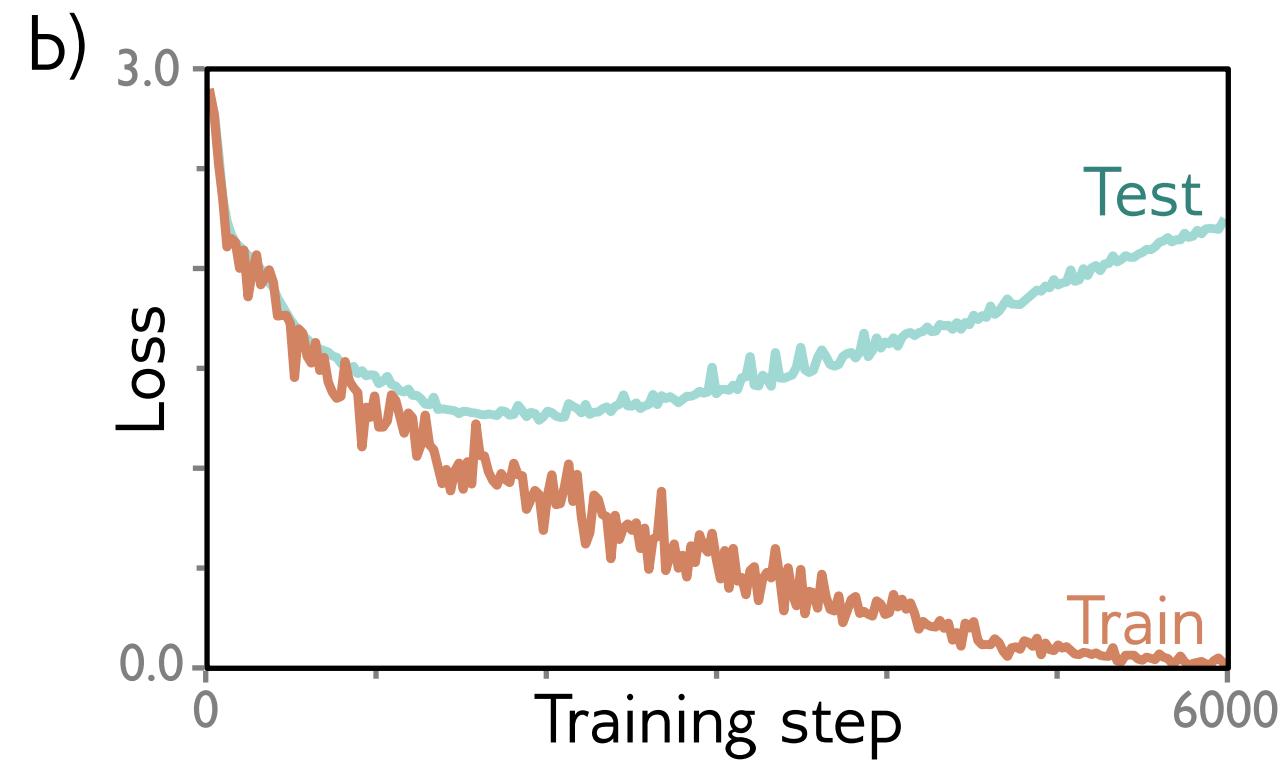
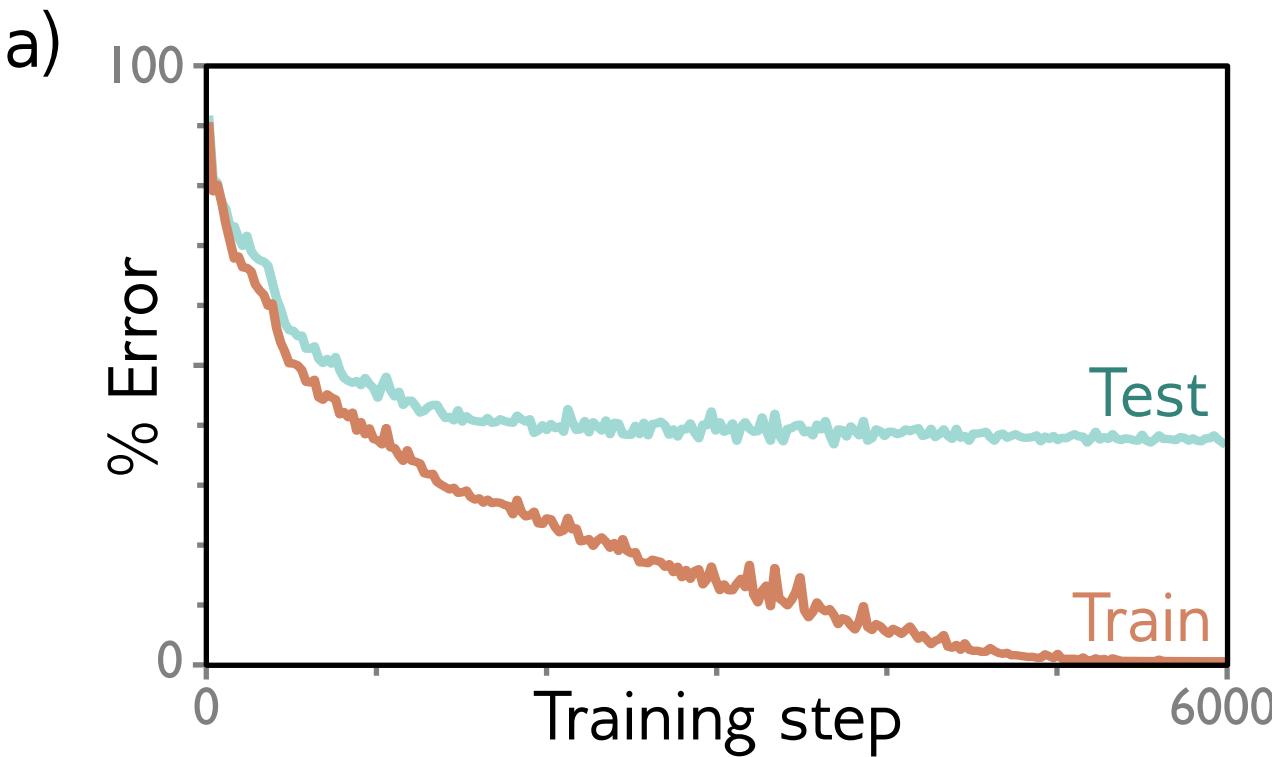
Network

- 40 inputs
- 10 outputs
- 4000 training examples (~400 training examples per class)
- Two hidden layers
 - 100 hidden units each
- SGD with batch size 100, learning rate 0.1
- 6000 steps (60 Epochs)

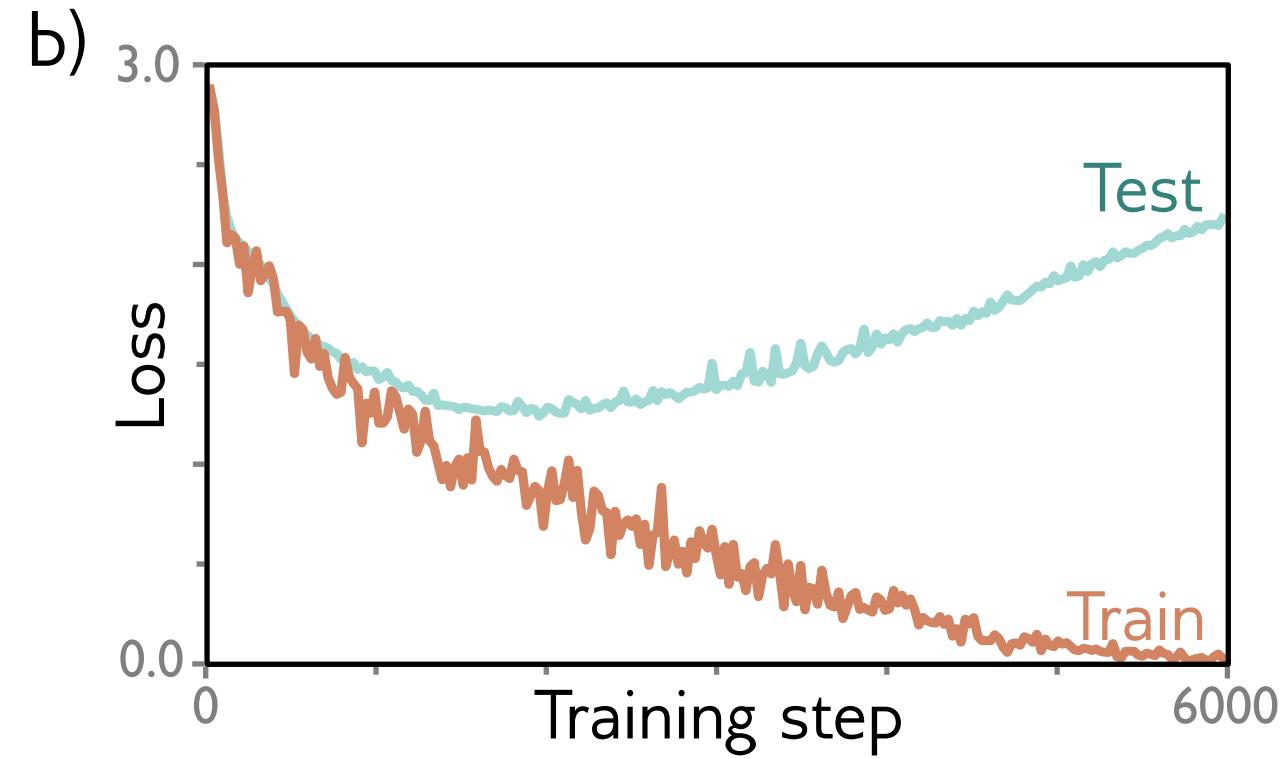
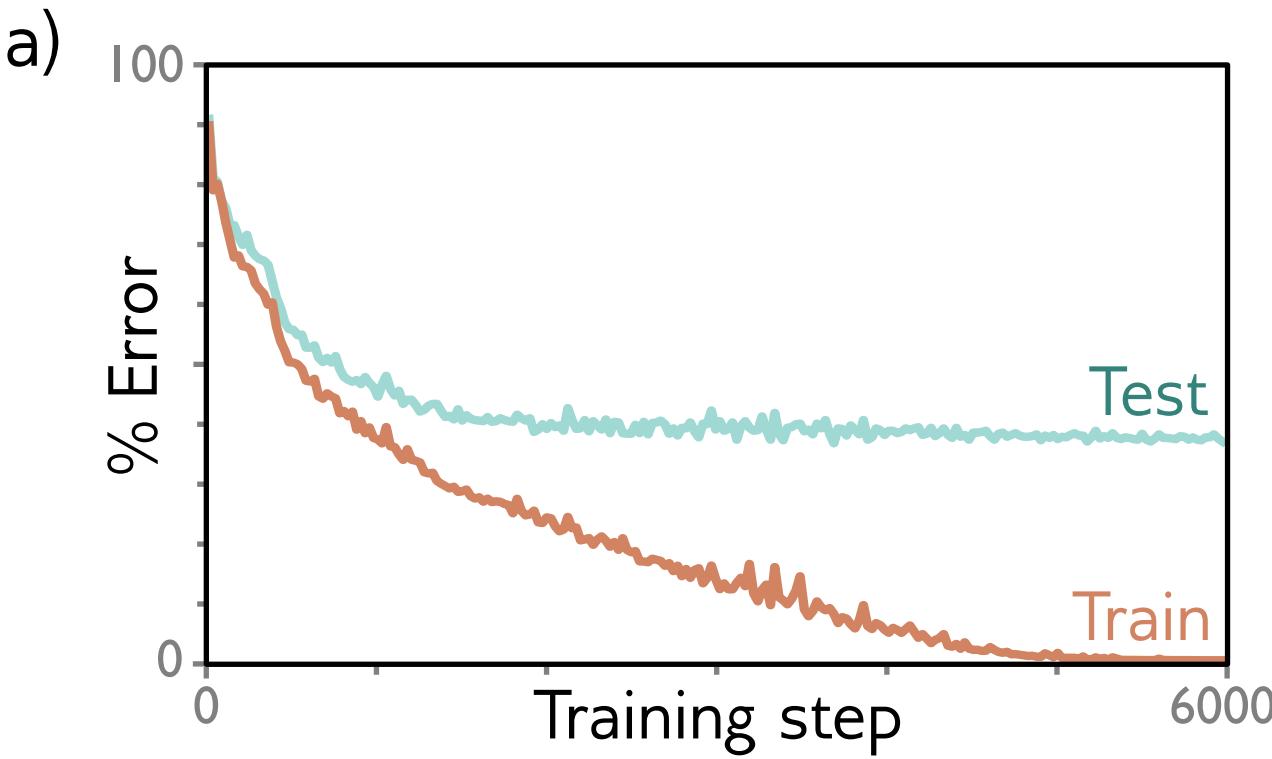
Results



Need to use separate test data



Need to use separate test data

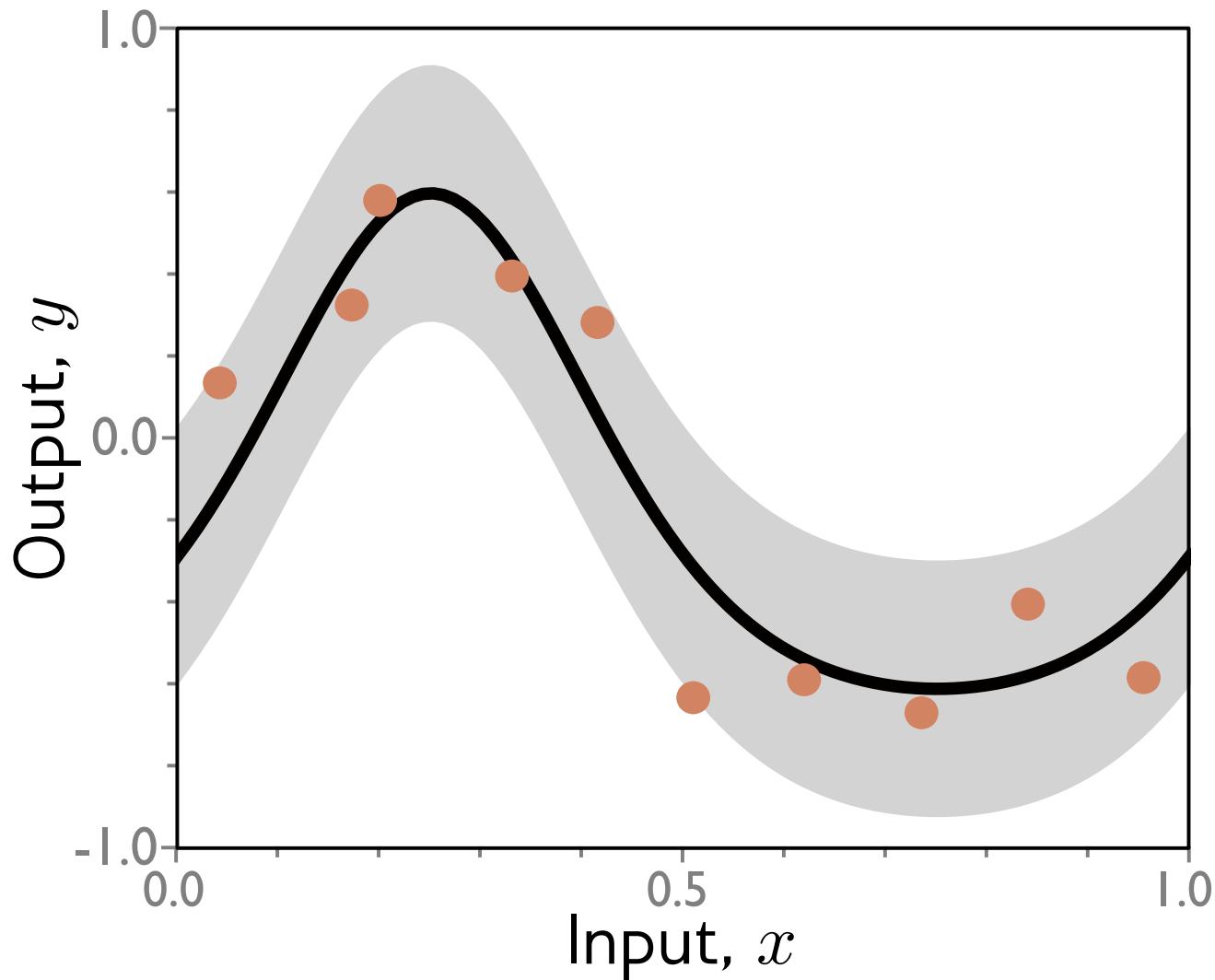


The model has not generalized well to the new data

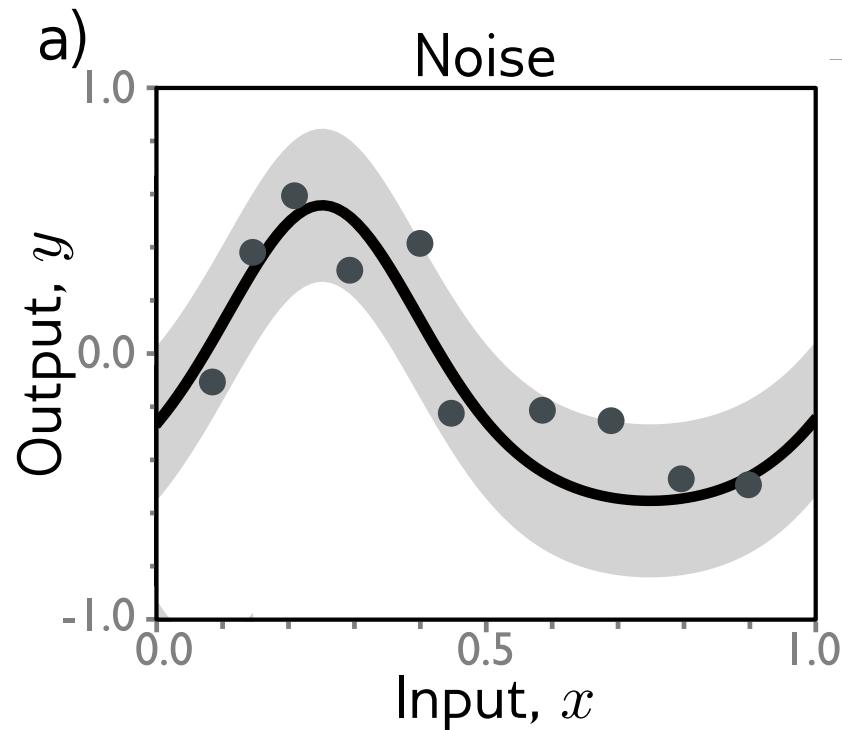
Measuring performance

- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Curse of dimensionality & weird properties of high dimensional space
- Choosing hyperparameters

Regression example



Noise, bias, and variance



- Noise in measurements
- Some variables not observed
- Data mislabeled

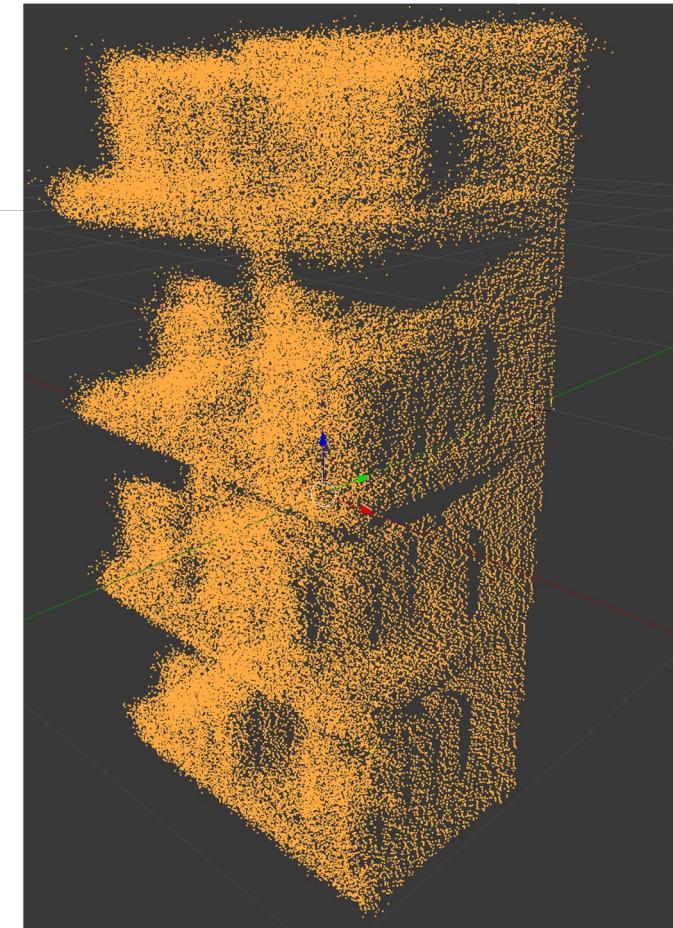
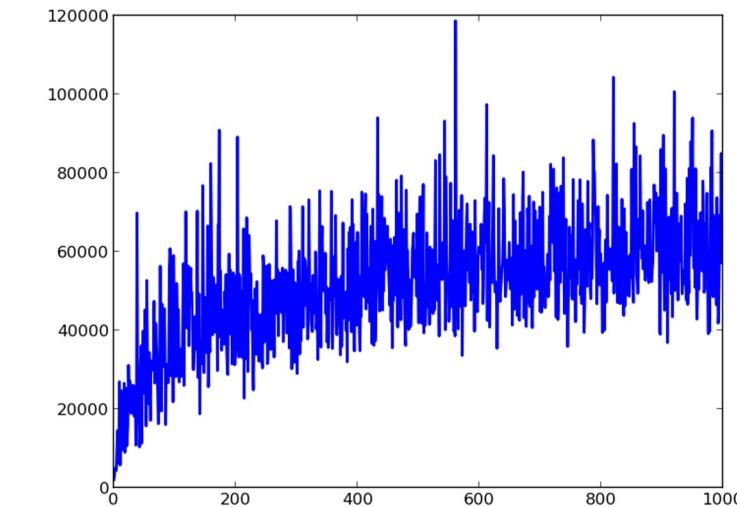
Noise is inherent uncertainty in the true mapping from input to output

Noise = Aleatoric Uncertainty

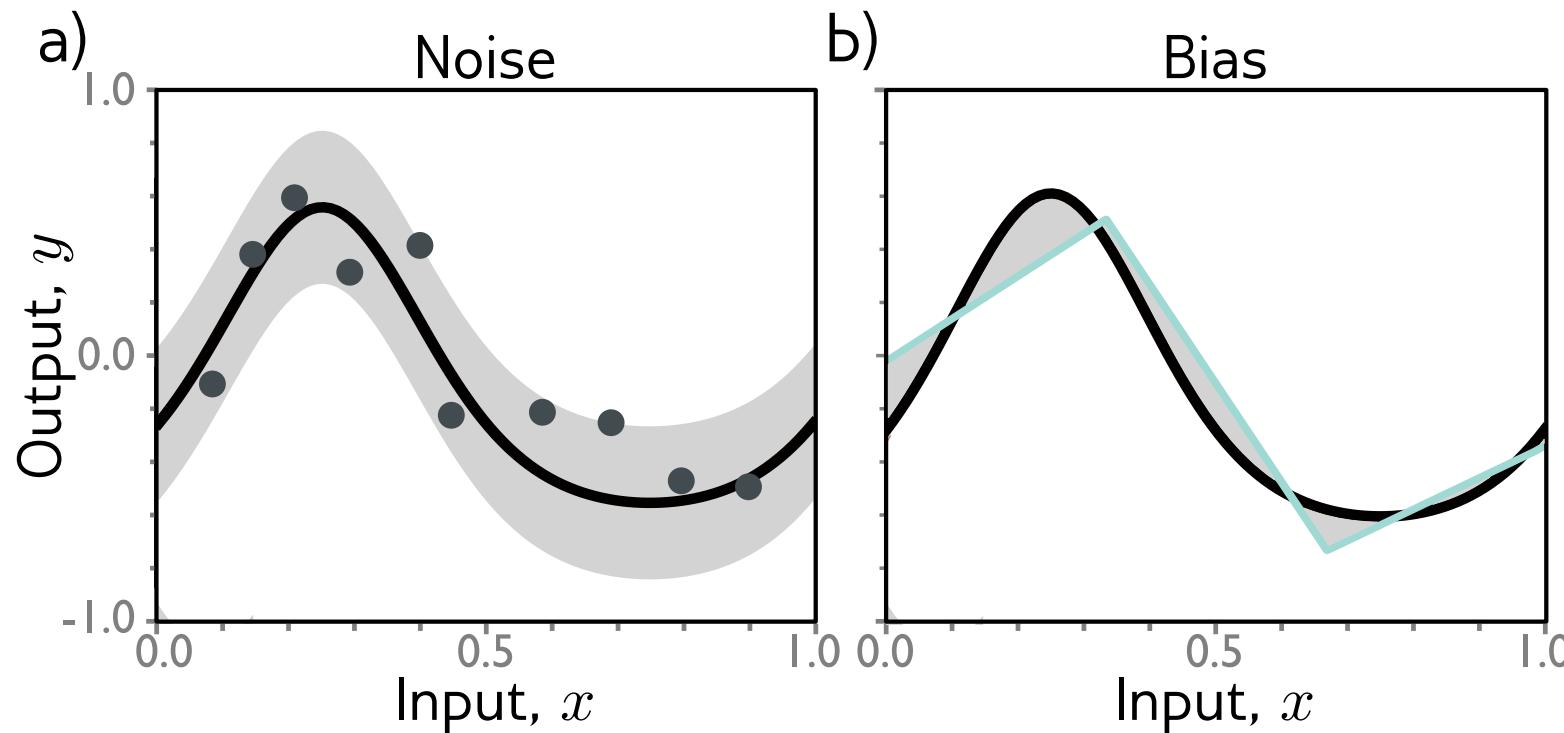
Noise that are inherent in the **observations**

Cannot be reduced with more data

Example: a **sensor noise** or motion noise, resulting in uncertainty which cannot be reduced **even if more data were to be collected.**

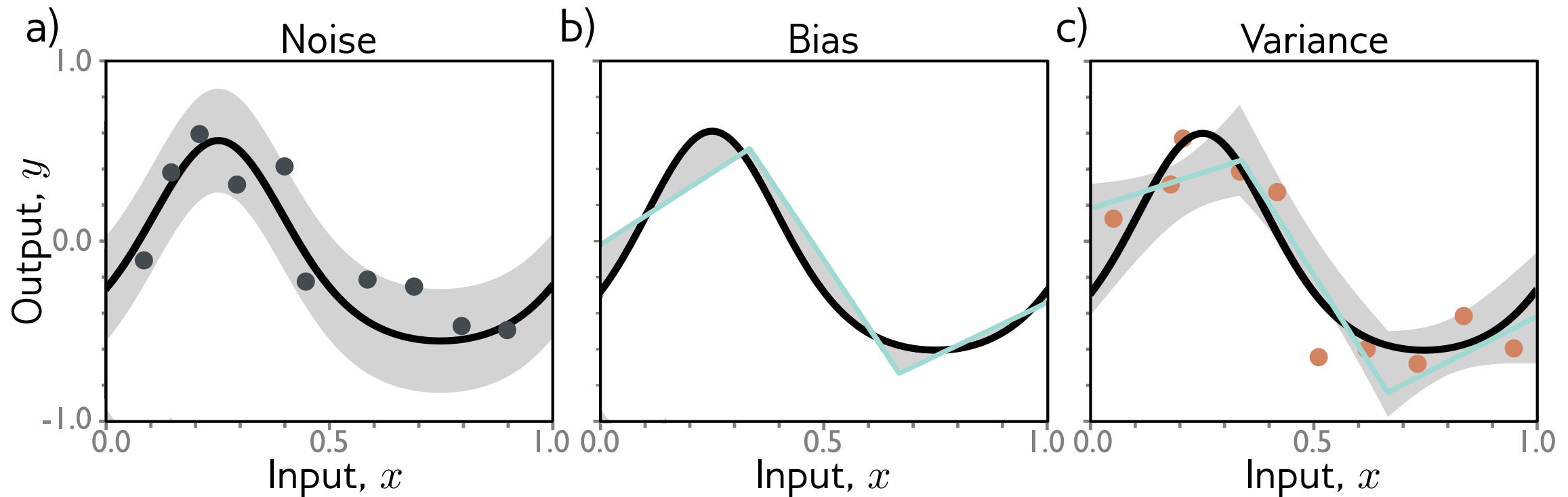


Noise, **bias**, and variance



Bias is systematic deviation from the mean of the function we are modeling due to limitations in our model

Noise, bias, and variance



Variance is the uncertainty in fitted model due to choice of training set

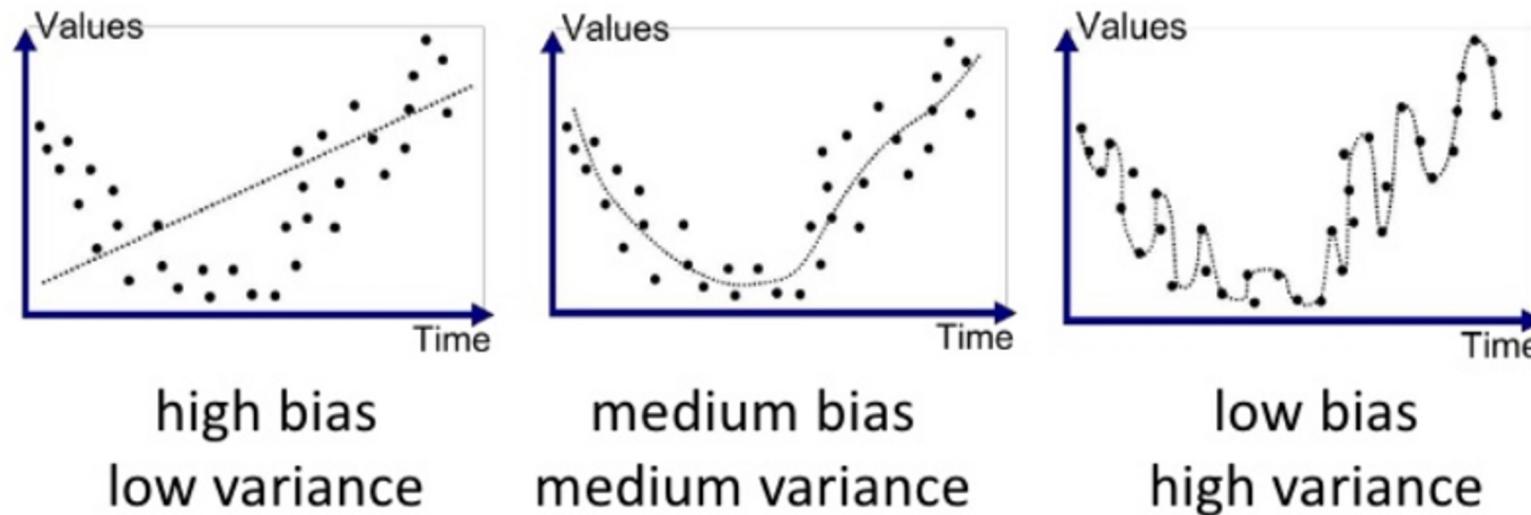
Bias + Variance = Epistemic Uncertainty

Uncertainty in the parameters of a model

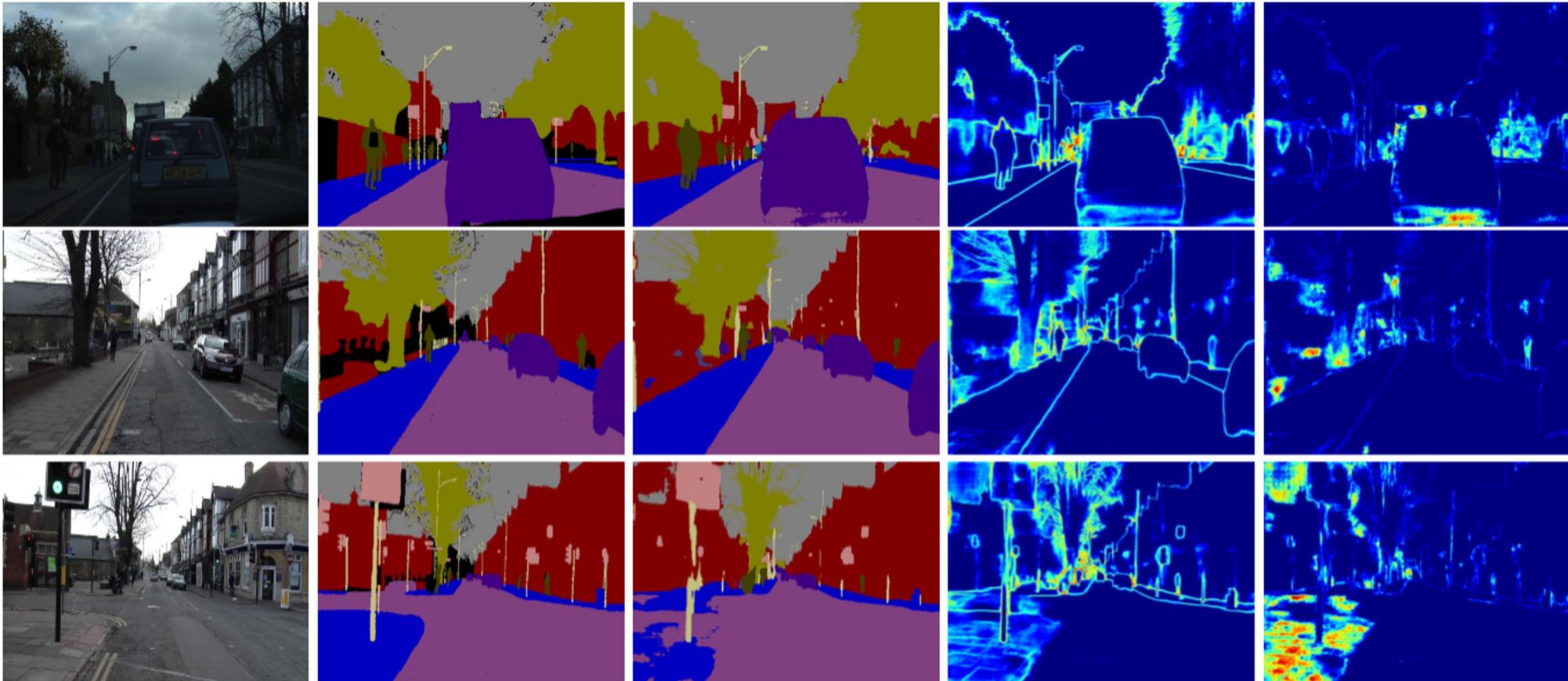
Captures our ignorance about which model generated our collected data

Expected to be high when the system encounters examples that differ from the training data (out-of-distribution)

Can (eventually) be reduced to zero given more data



Epistemic and Aleatoric Uncertainty



(a) Input Image

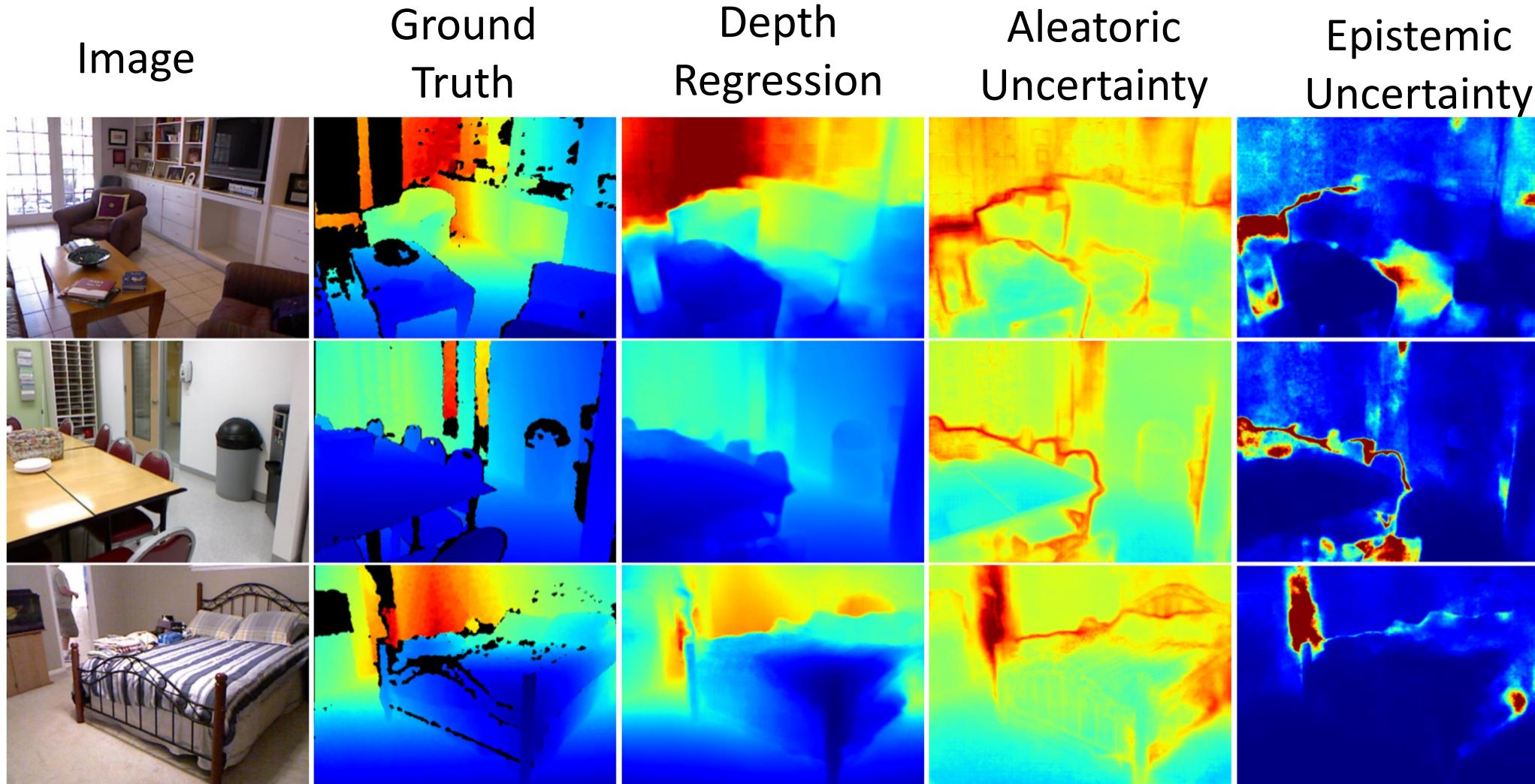
(b) Ground Truth

(c) Semantic
Segmentation

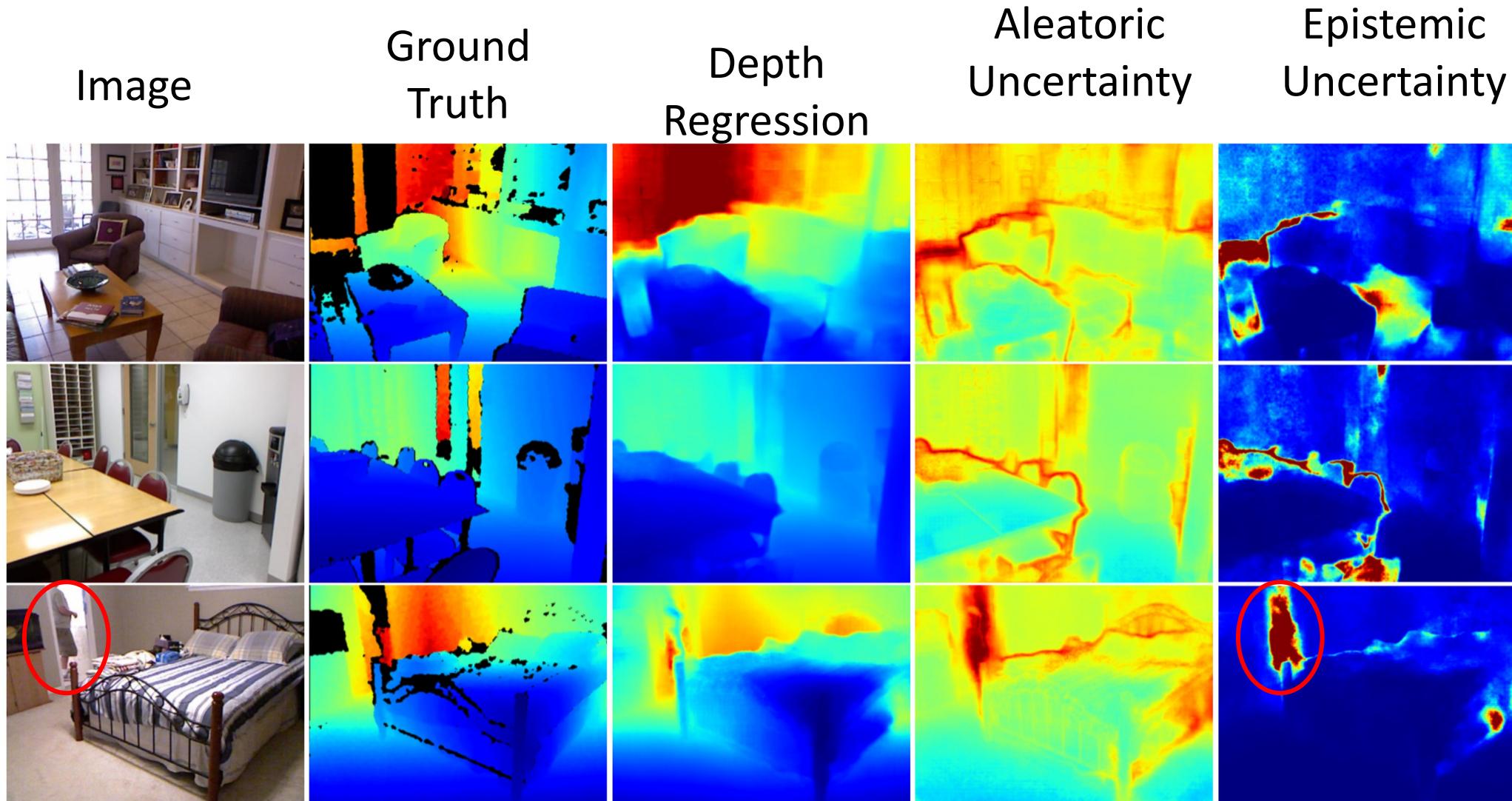
(d) Aleatoric
Uncertainty

(e) Epistemic
Uncertainty

Epistemic and Aleatoric Uncertainty



Epistemic and Aleatoric Uncertainty

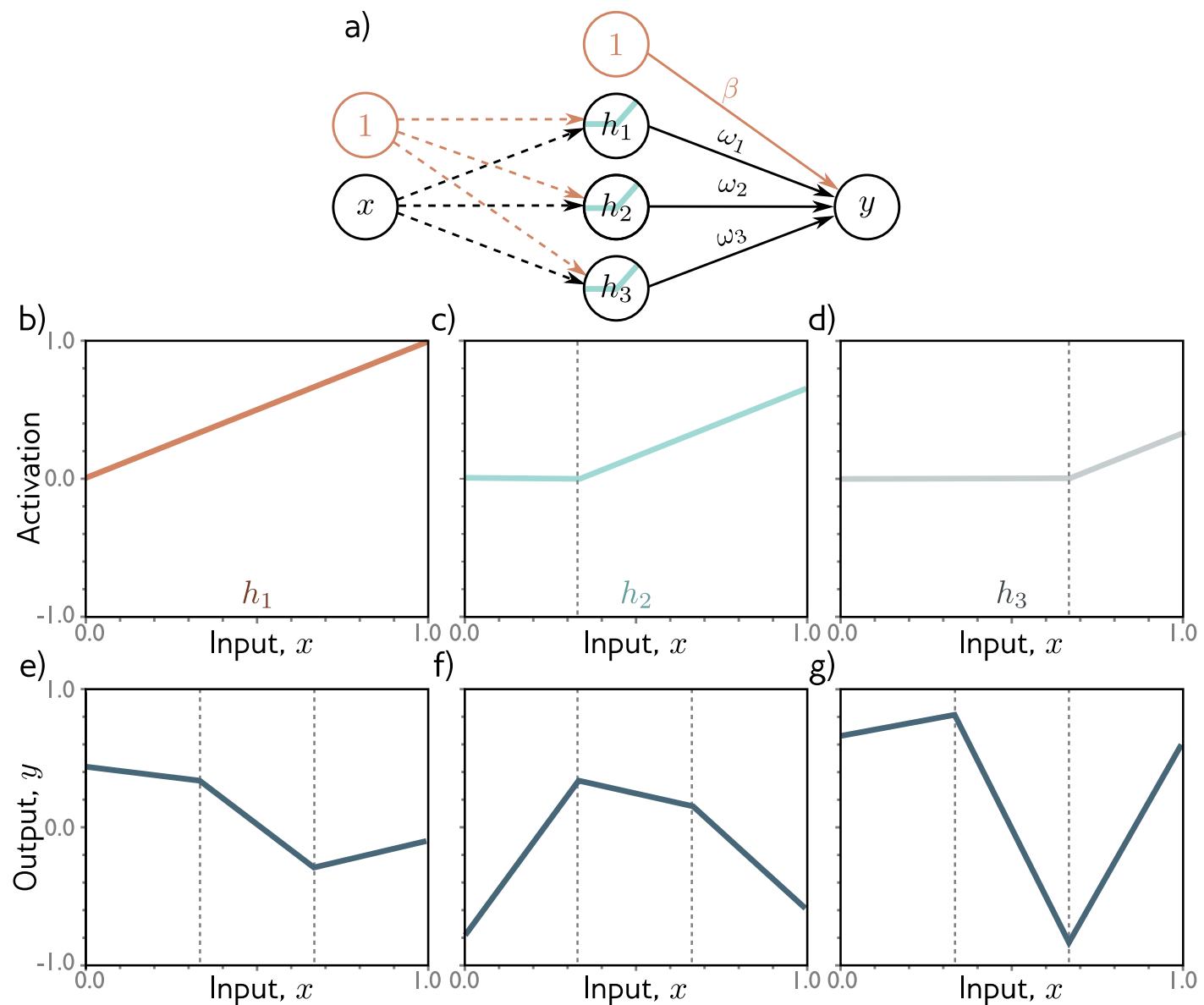


Noise, bias, and variance

- Variance is the uncertainty in fitted model due to choice of training set
- Bias is systematic deviation from the mean of the function we are modeling due to limitations in our model
- Noise is inherent uncertainty in the true mapping from input to output

Toy model

- K hidden units
- First layer fixed so “joints” divide interval evenly
- Second layer trained
- But... now linear in h
 - so convex cost function
 - can find best solution in closed-form



Least squares regression only

$$L[x] = (f[x, \phi] - y[x])^2$$

- We can show that:

$$\mathbb{E}_{\mathcal{D}} [\mathbb{E}_y [L[x]]] = \underbrace{\mathbb{E}_{\mathcal{D}} [(f[x, \phi[\mathcal{D}]] - f_\mu[x])^2]}_{\text{variance}} + \underbrace{(f_\mu[x] - \mu[x])^2}_{\text{bias}} + \underbrace{\sigma^2}_{\text{noise}}$$

Expectation over noise in training data

Expectation over noise in test data

Actual model

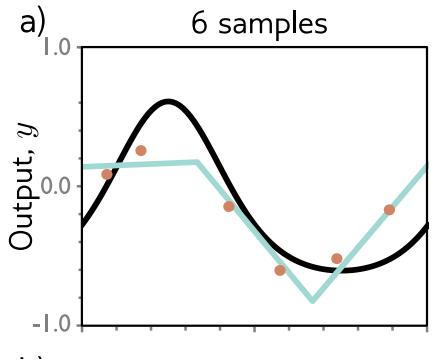
Best possible model if we had infinite data

True function

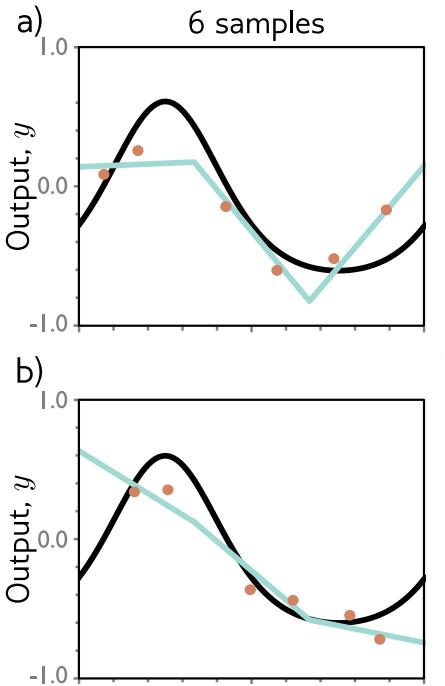
Measuring performance

- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Curse of dimensionality & weird properties of high dimensional space
- Choosing hyperparameters

Variance

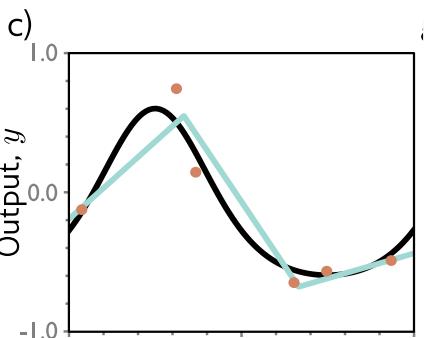
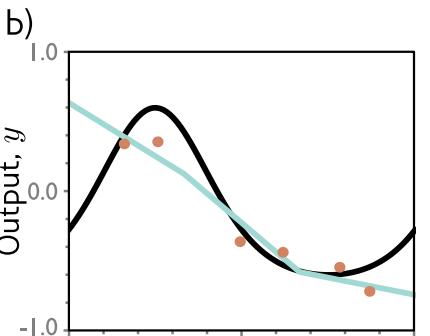
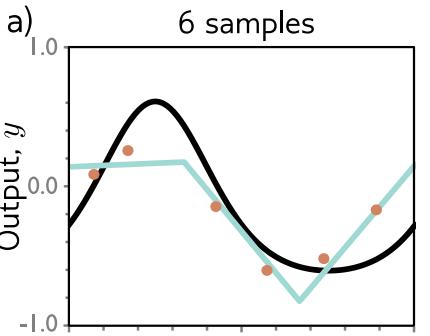


Variance



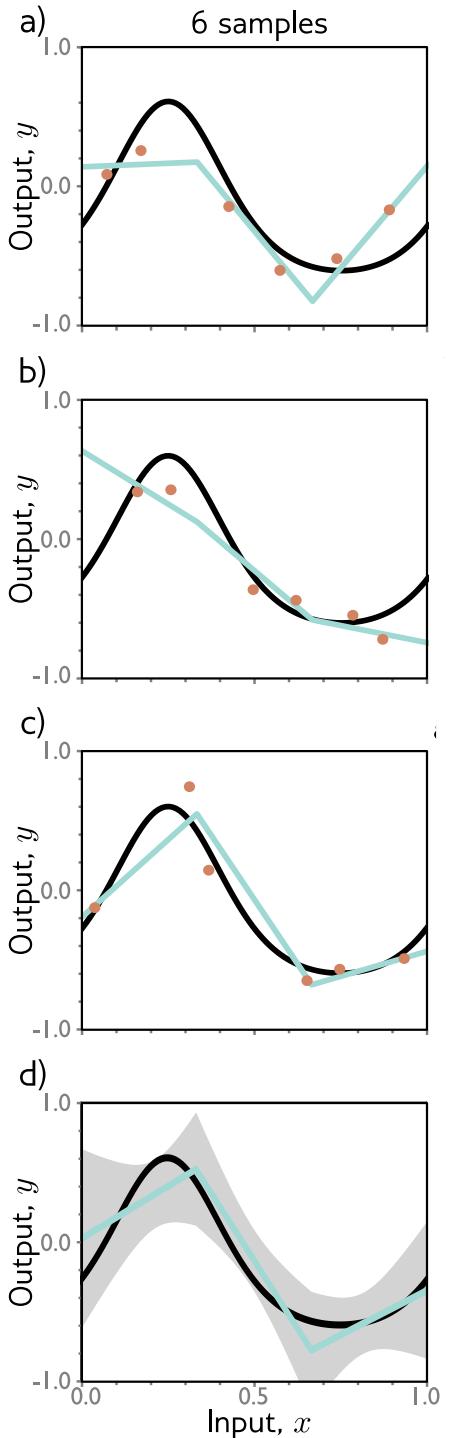
Draw sample 2

Variance



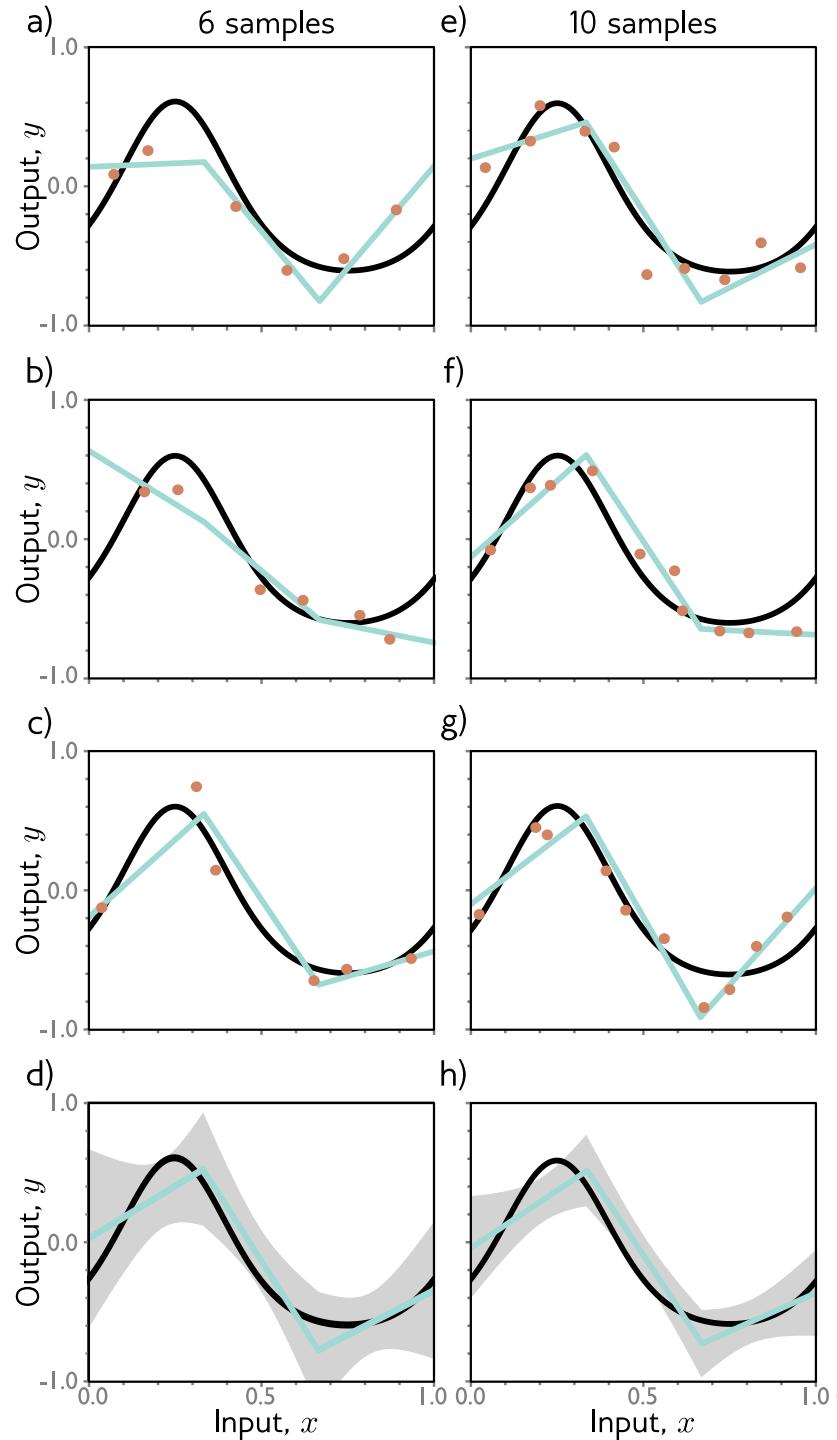
Draw sample 2

Variance

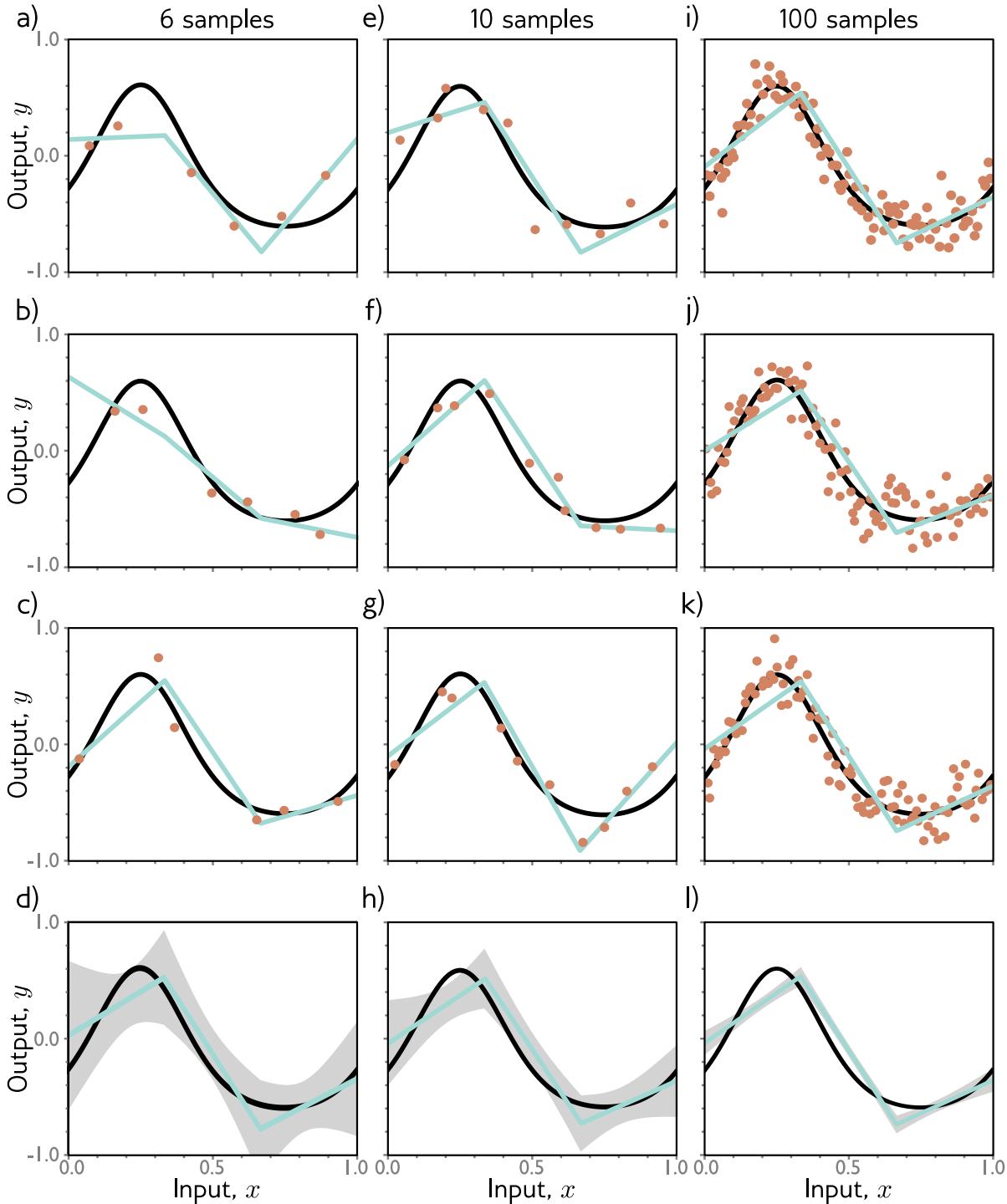


Variance over 3
draws

Variance



Variance

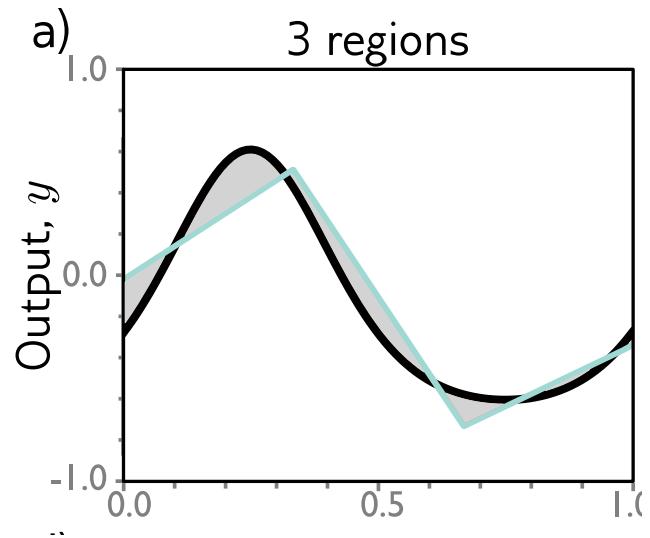


Can reduce
variance by
adding more
samples

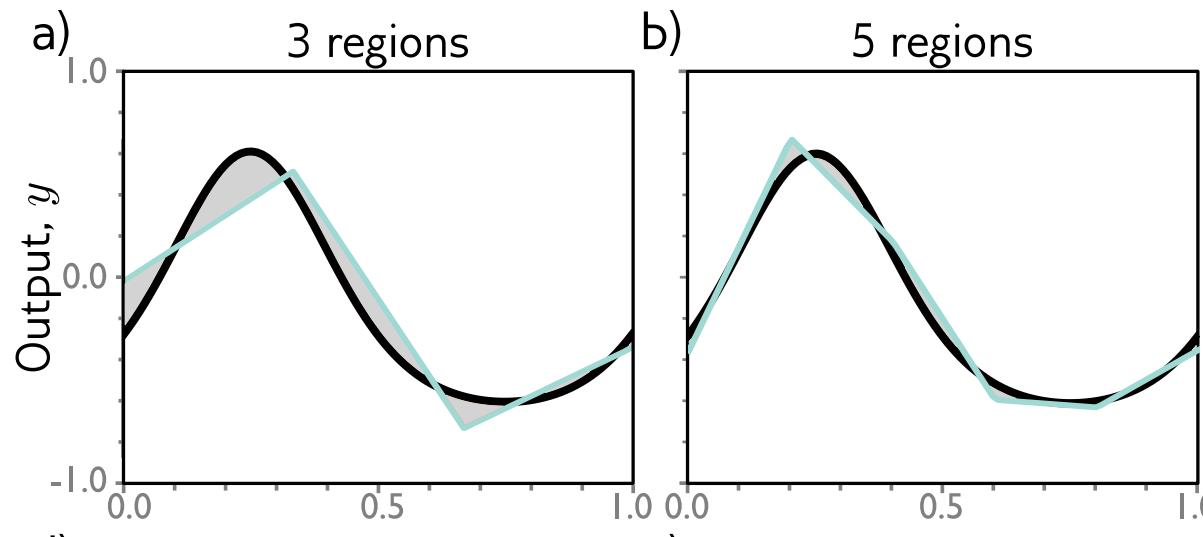
Measuring performance

- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Curse of dimensionality & weird properties of high dimensional space
- Choosing hyperparameters

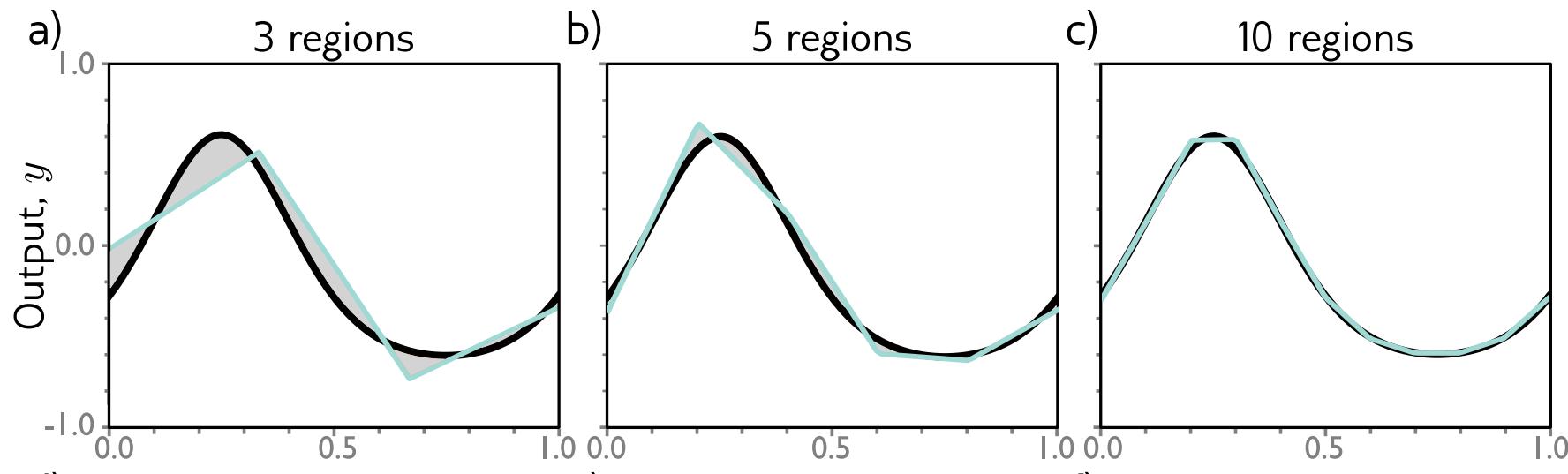
Reducing bias



Reducing bias

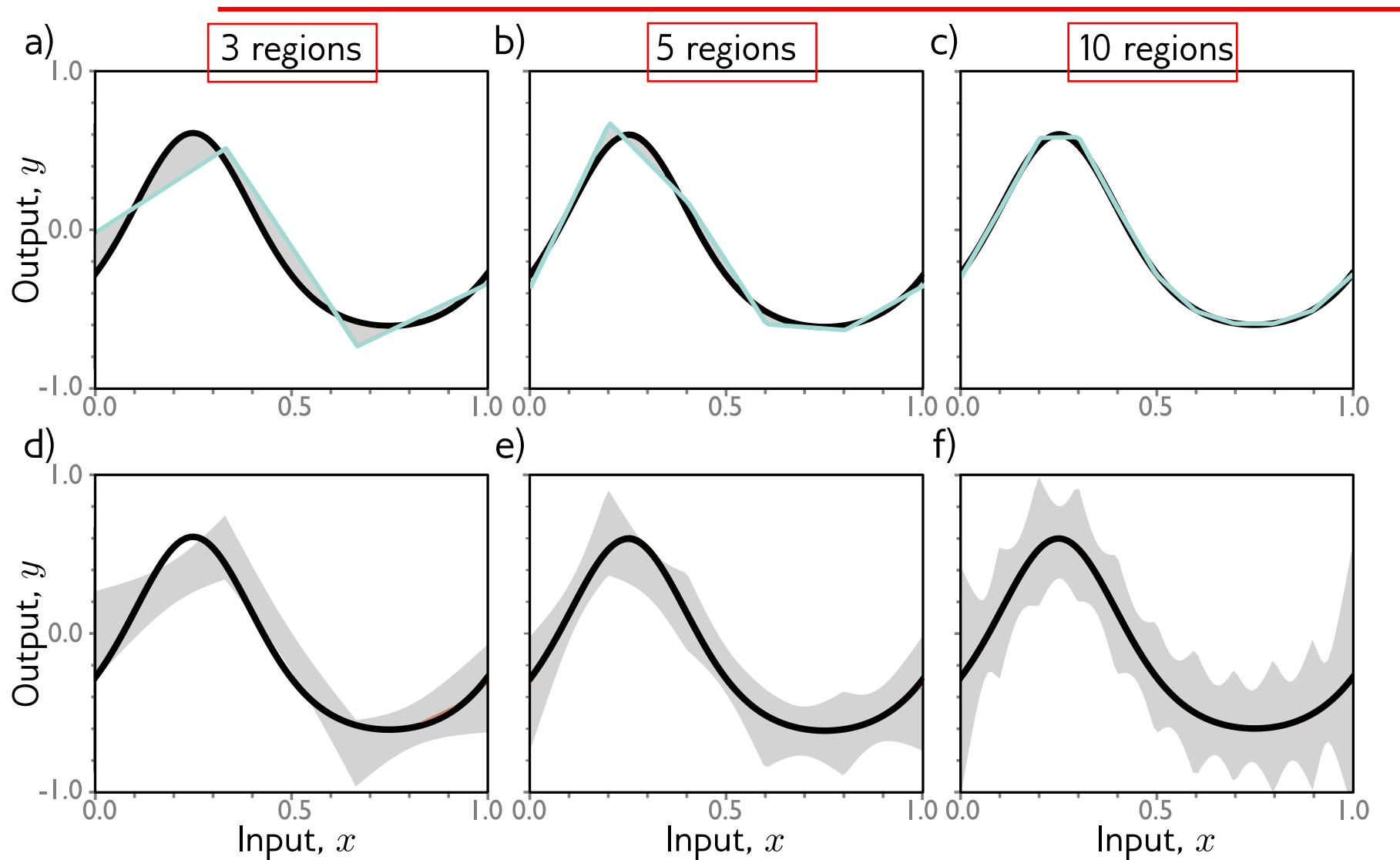


Reducing bias



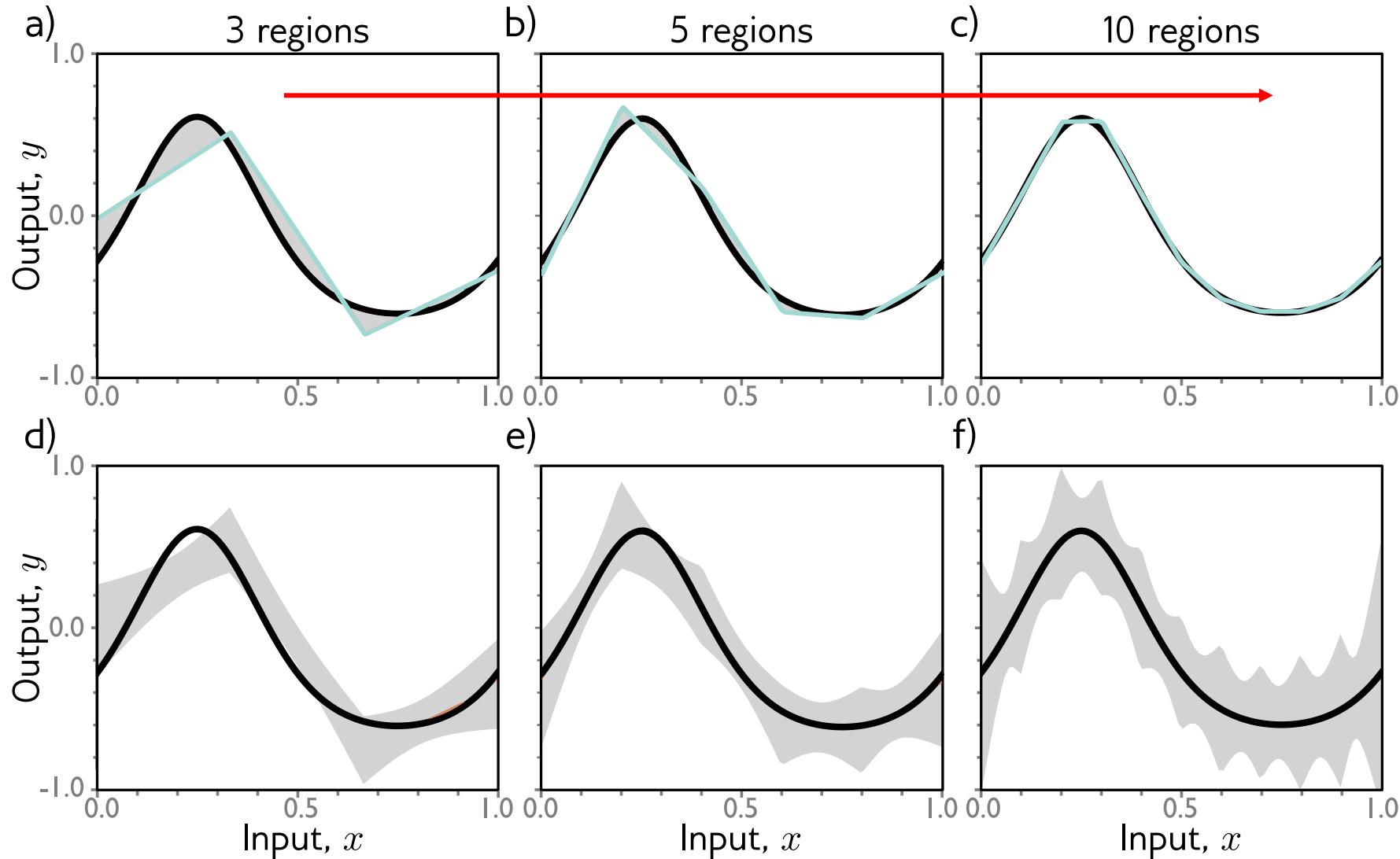
Reducing bias

Increasing Capacity



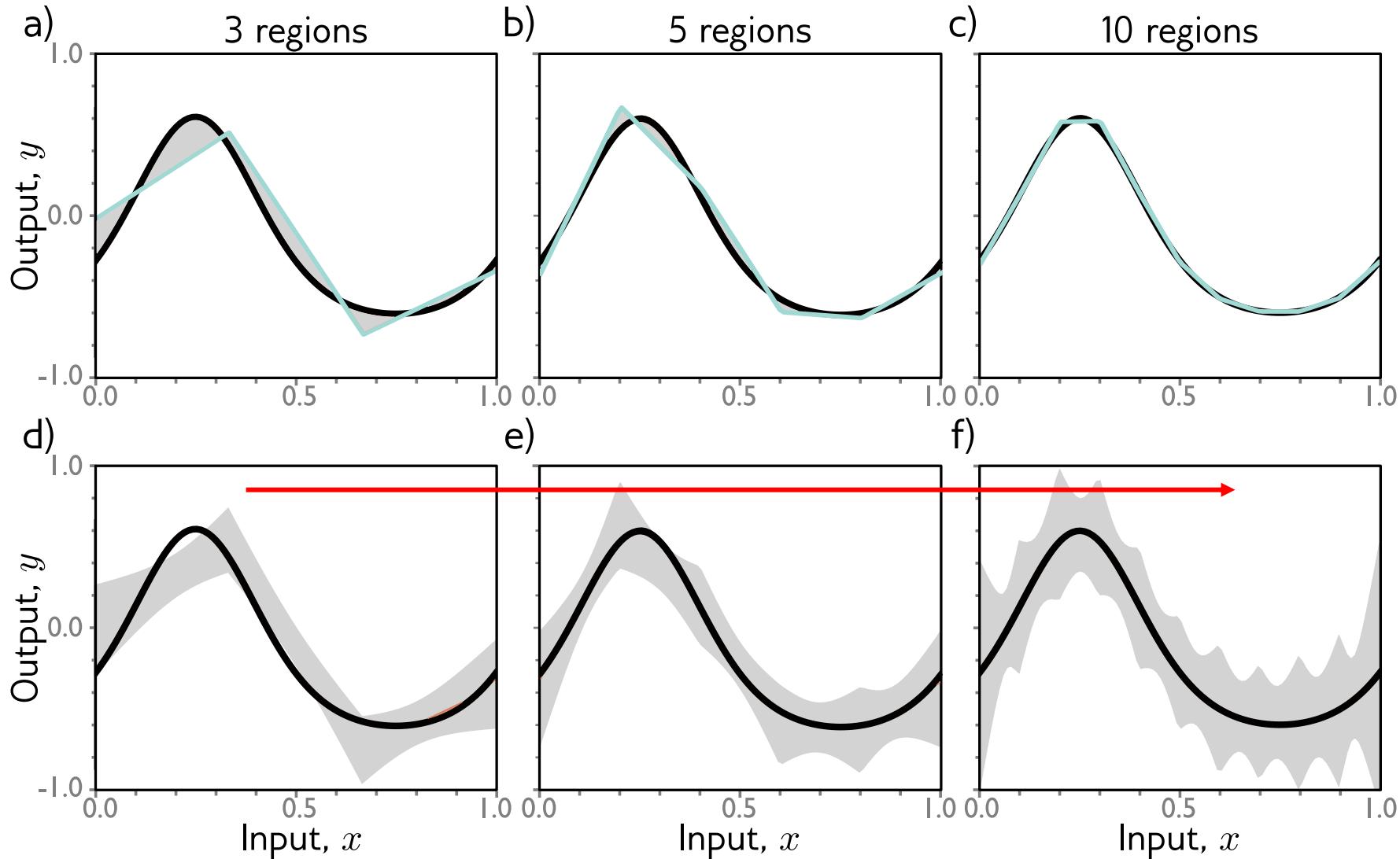
Reducing bias

Decreasing bias

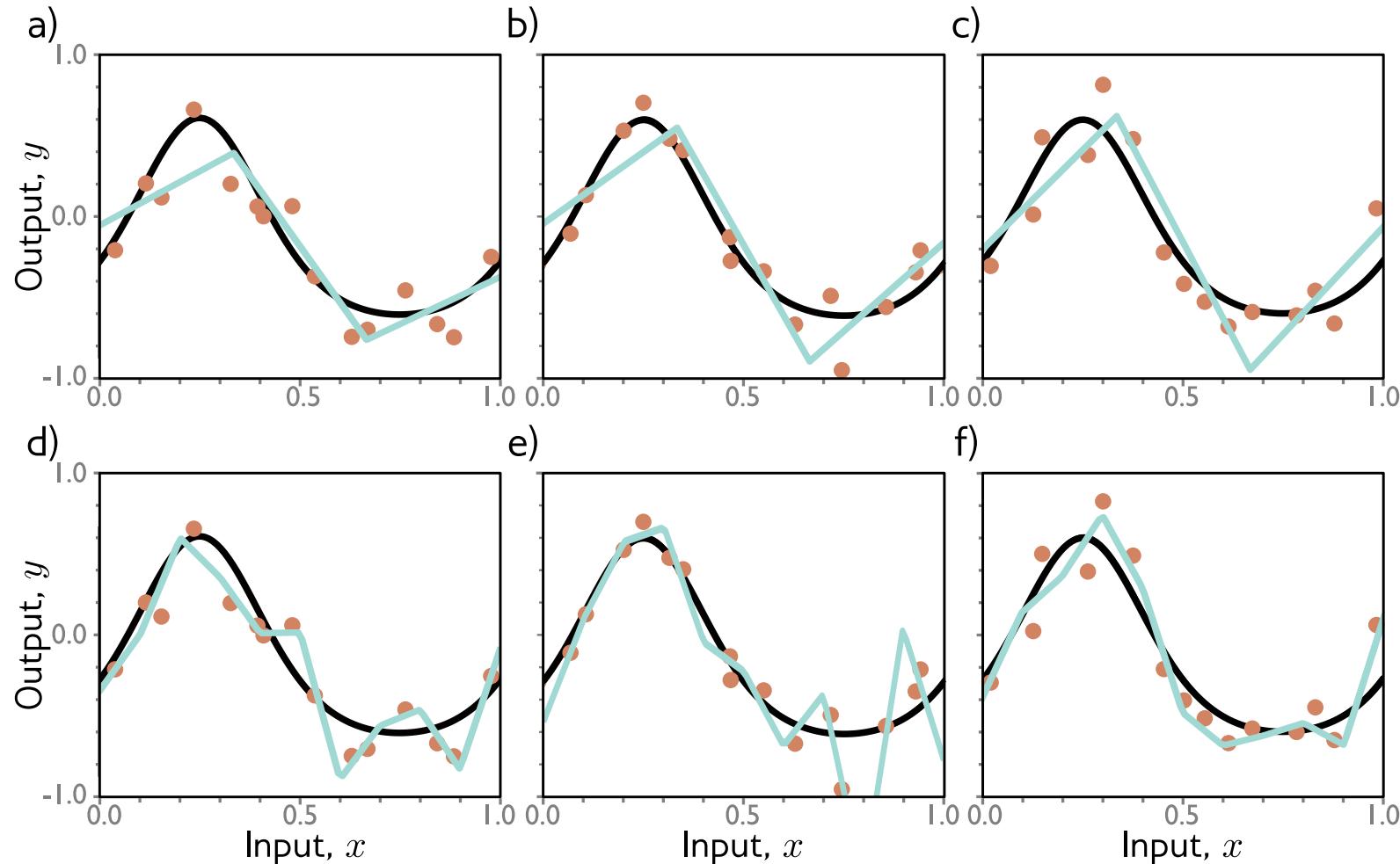


Reducing bias

Increasing Variance

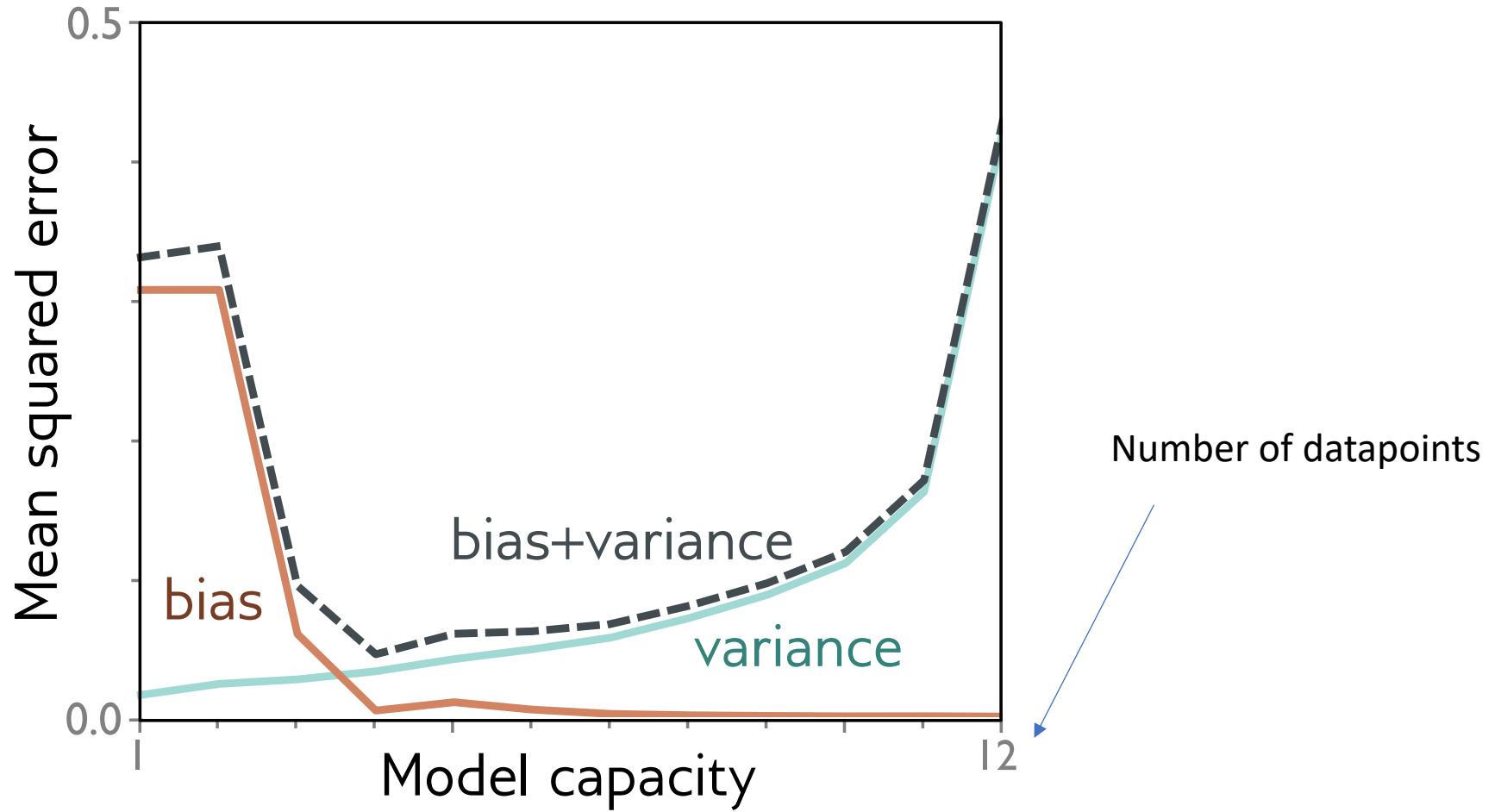


Why does variance increase? Overfitting



Describes the training data better, but not the true underlying function (black curve)

Bias and variance trade-off



Measuring performance

- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Curse of dimensionality & weird properties of high dimensional space
- Choosing hyperparameters

Questions?

Double Descent

“The discovery of double descent is recent, unexpected, and somewhat puzzling” – Prince, p.129

Double Descent

Interaction of two phenomena:

- 1) Test performance is worse and capacity equals data set size and can memorize exactly (interpolation threshold)

Double Descent

Interaction of two phenomena:

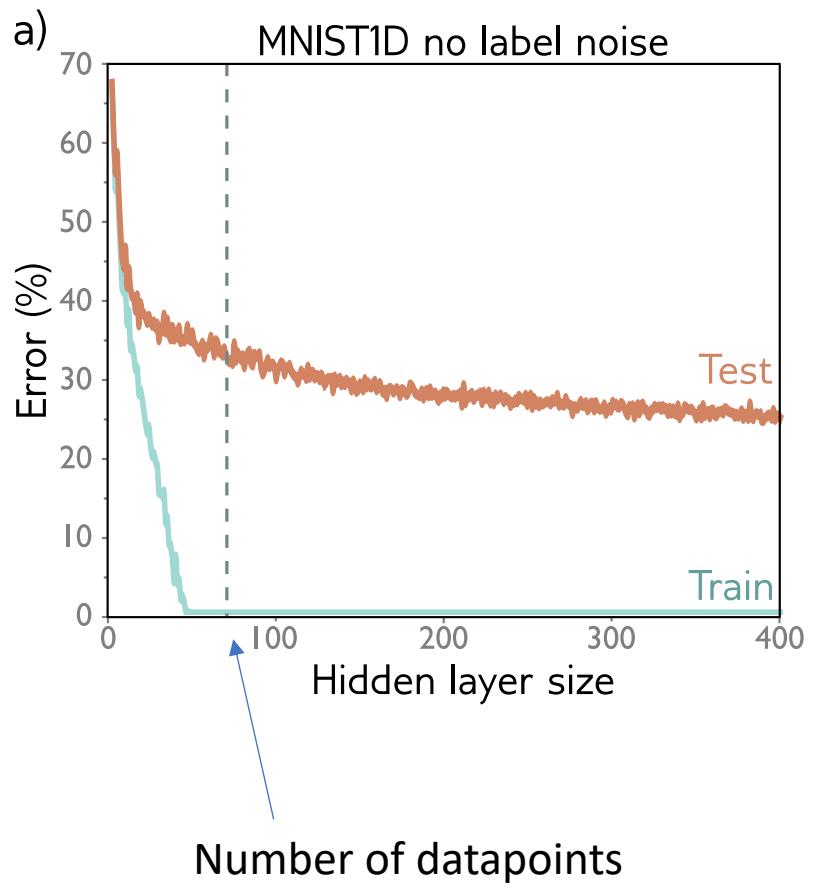
- 1) Test performance is worse and capacity equals data set size and can memorize exactly (interpolation threshold)
- 2) Test performance **improves** as we increase capacity beyond this point

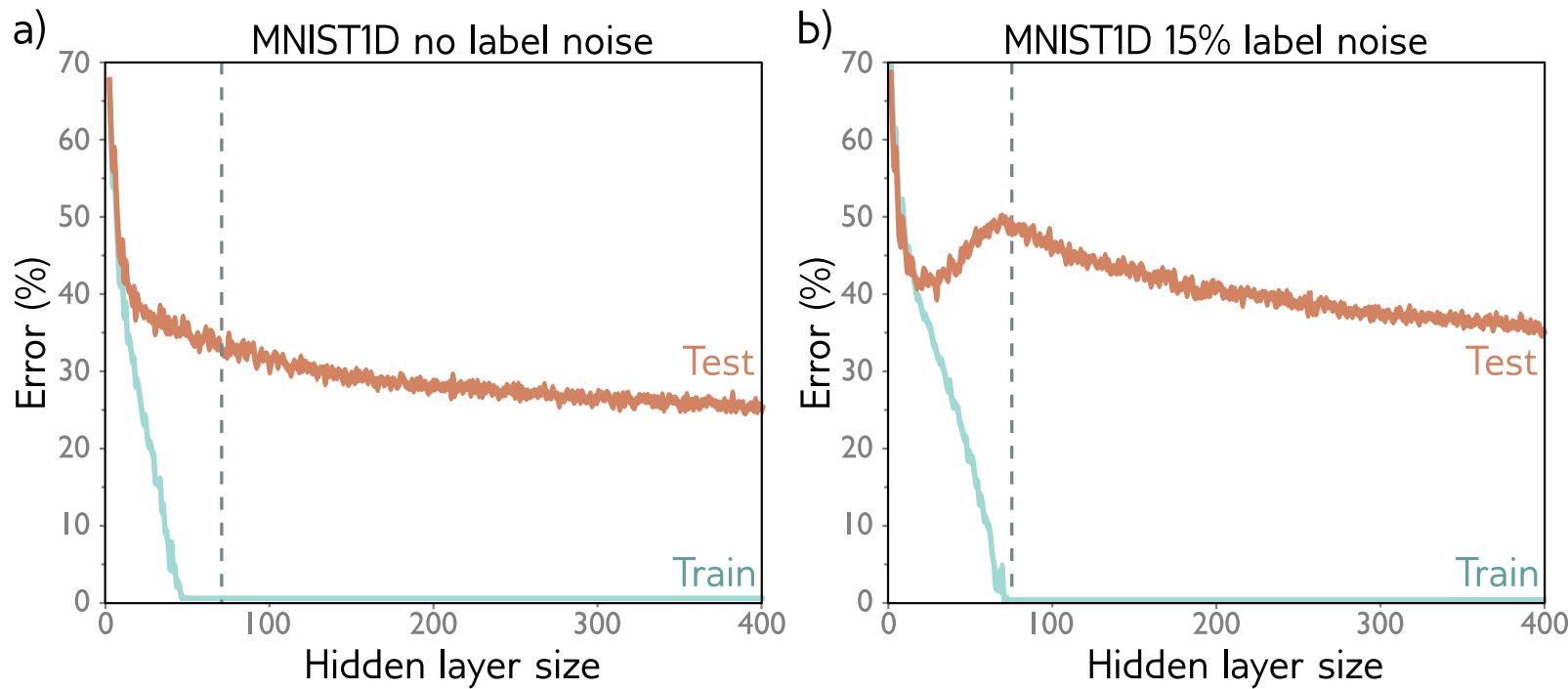
Double Descent

Interaction of two phenomena:

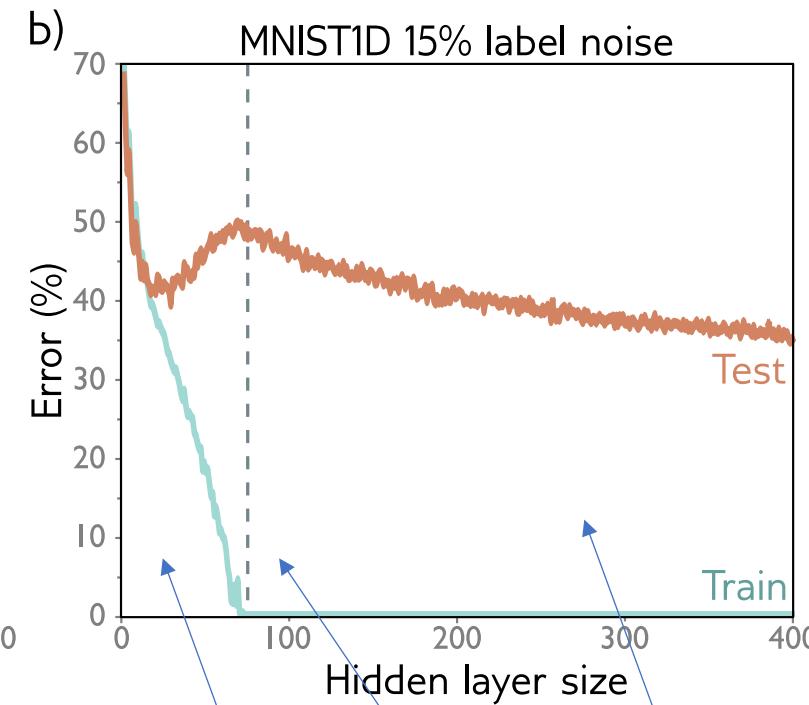
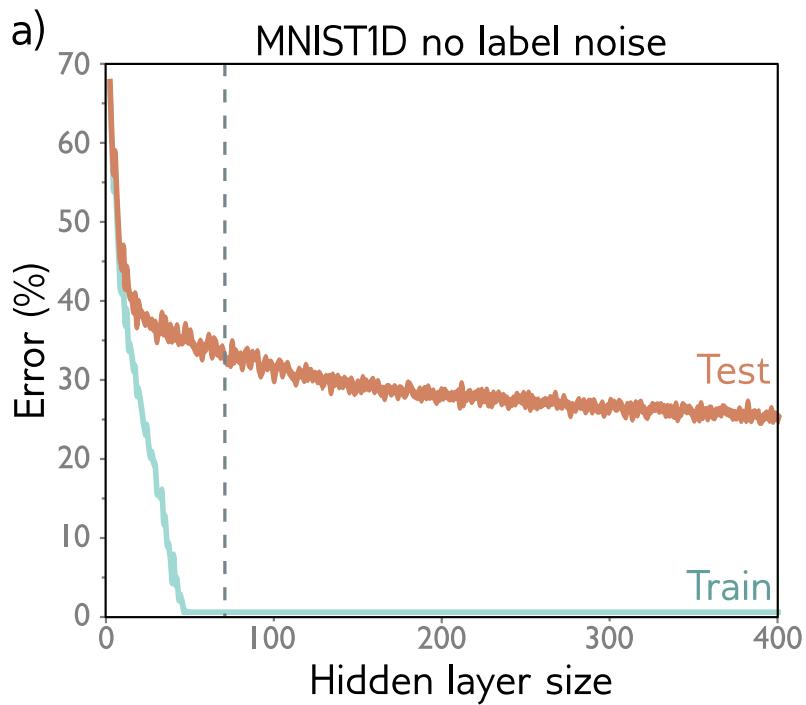
- 1) Test performance is worse and capacity equals data set size and can memorize exactly (interpolation threshold)
- 2) Test performance **improves** as we increase capacity beyond this point

Contradicts theory of bias-variance trade-off





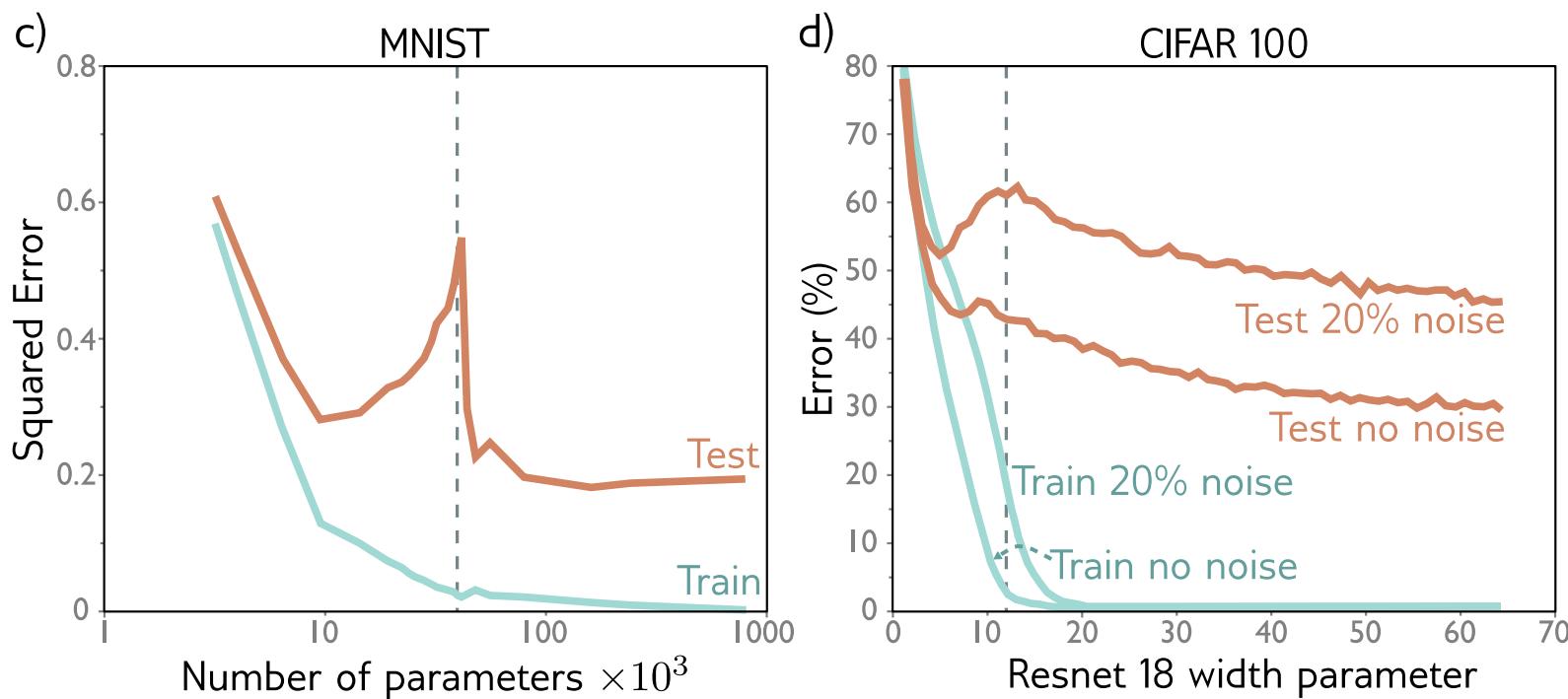
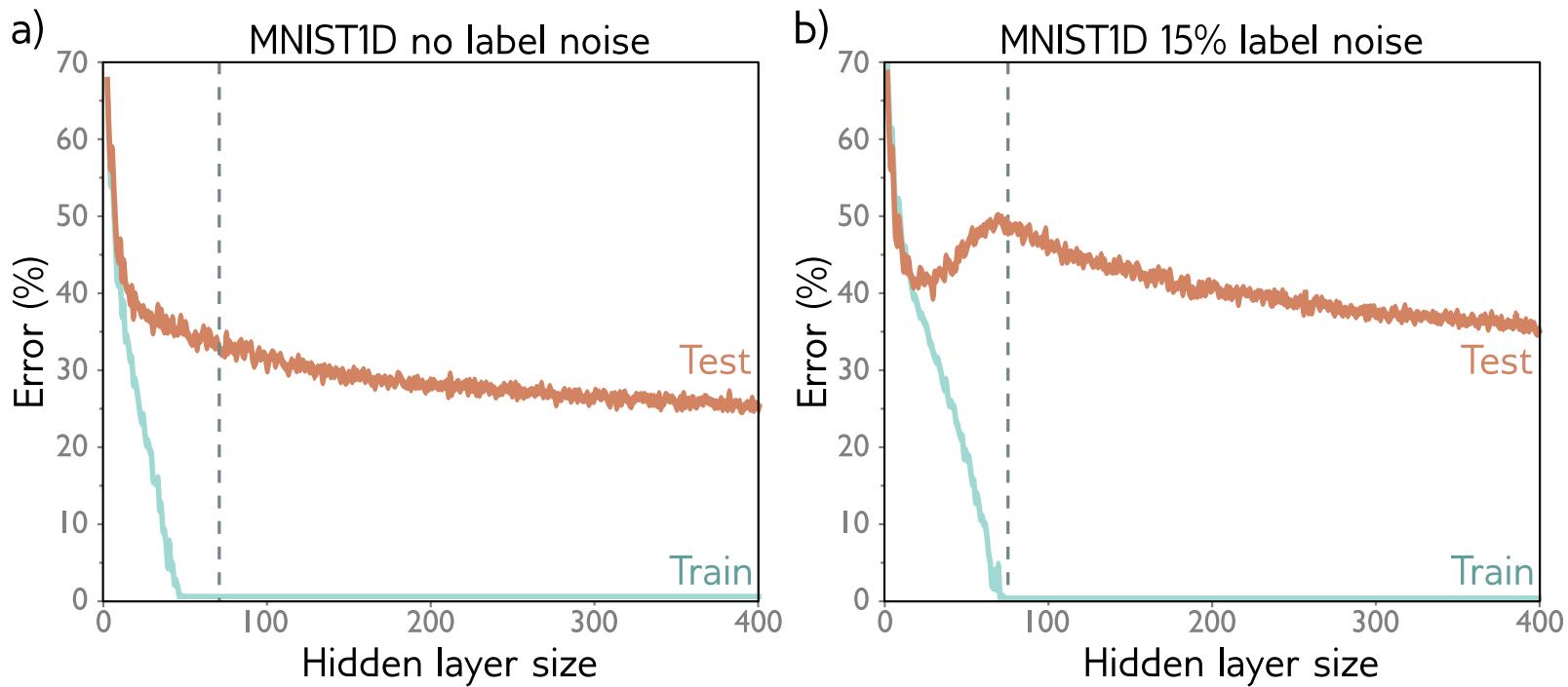
Double descent

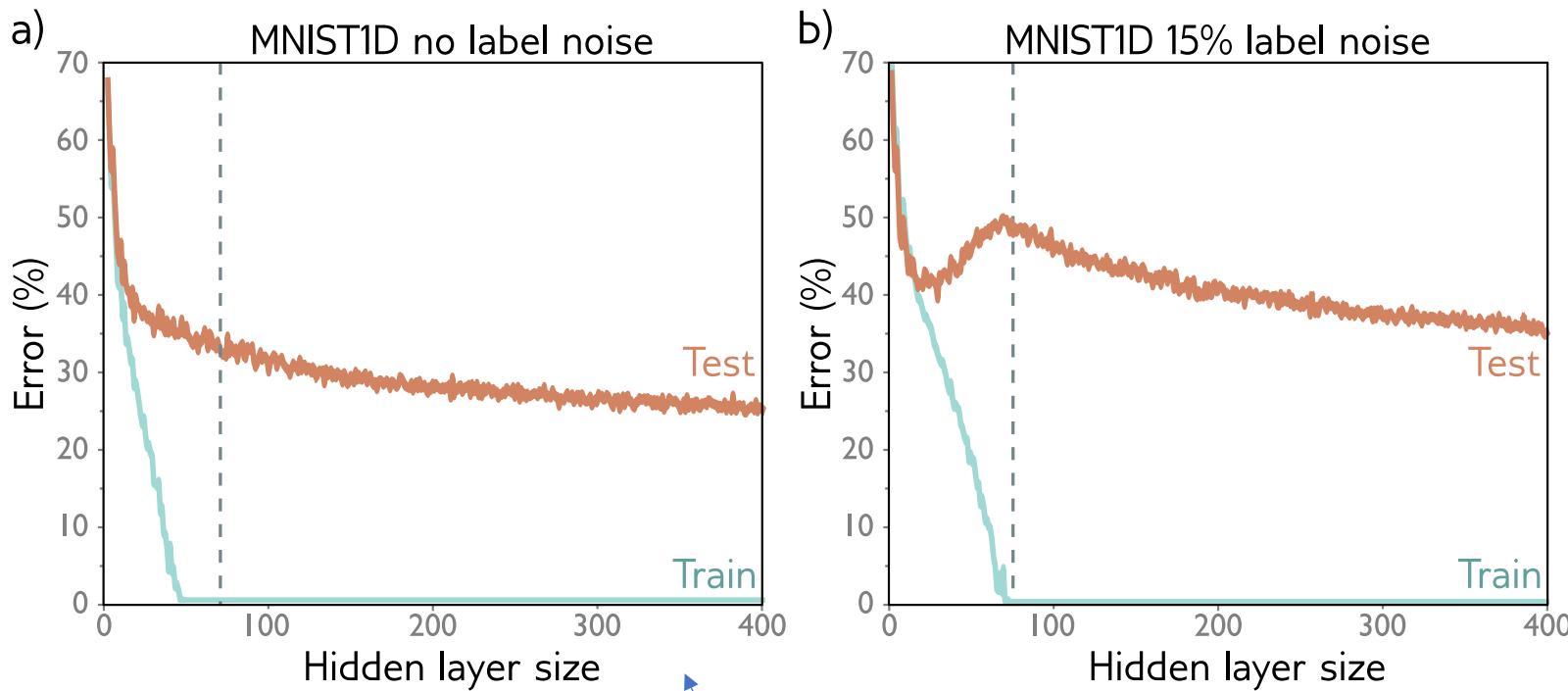


Classical or under-parameterized regime

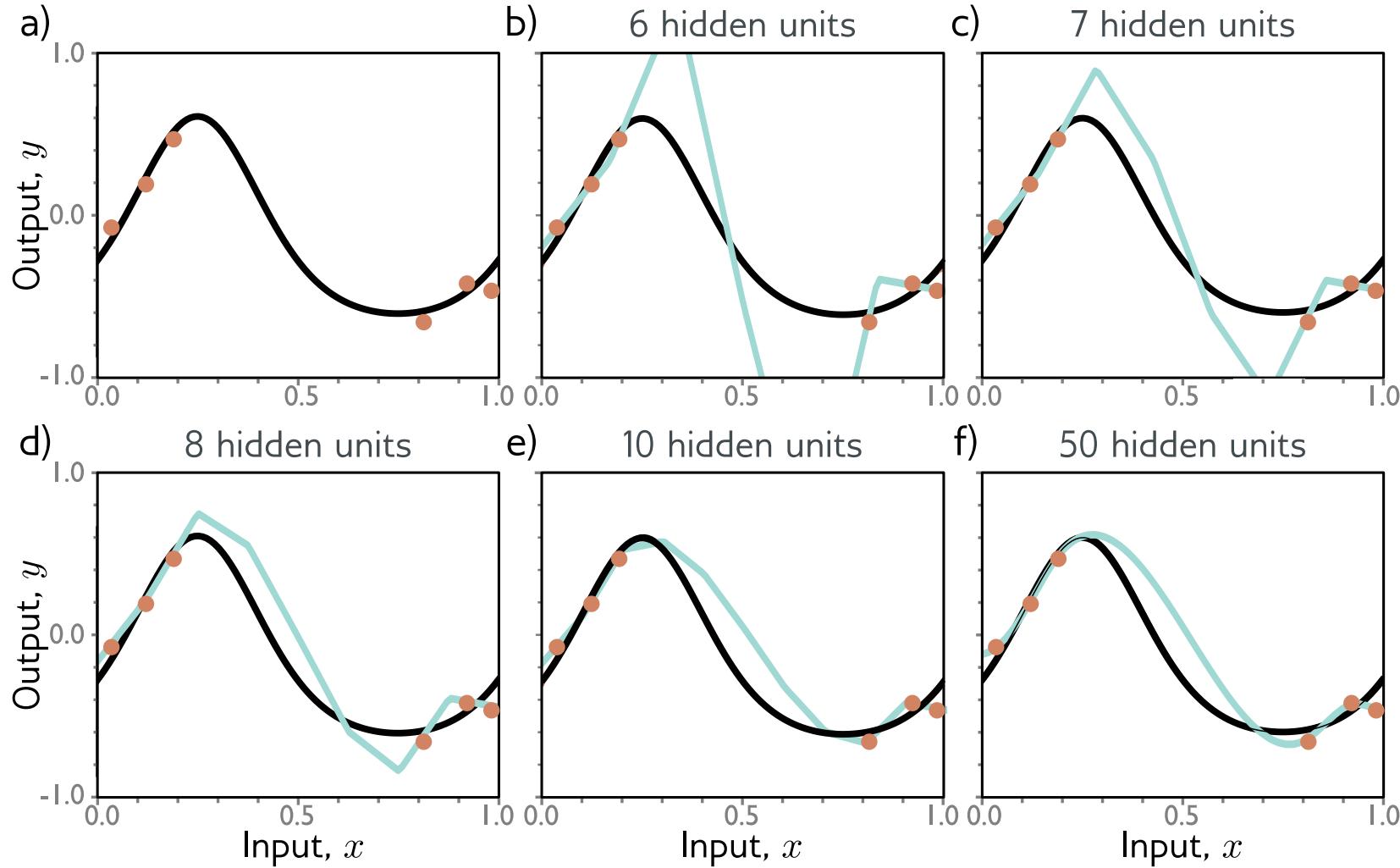
Modern or over-parameterized regime

Critical regime





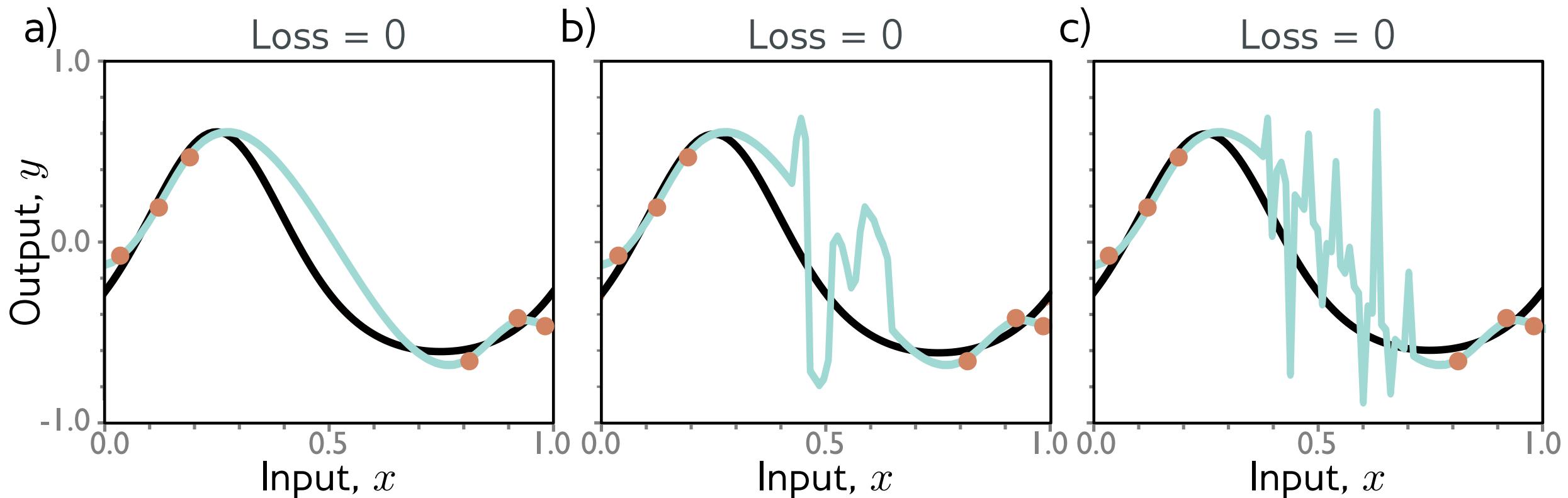
- Note that train data is very close to zero.
- Whatever is happening isn't happening at training data points
- Must be happening between the data points??



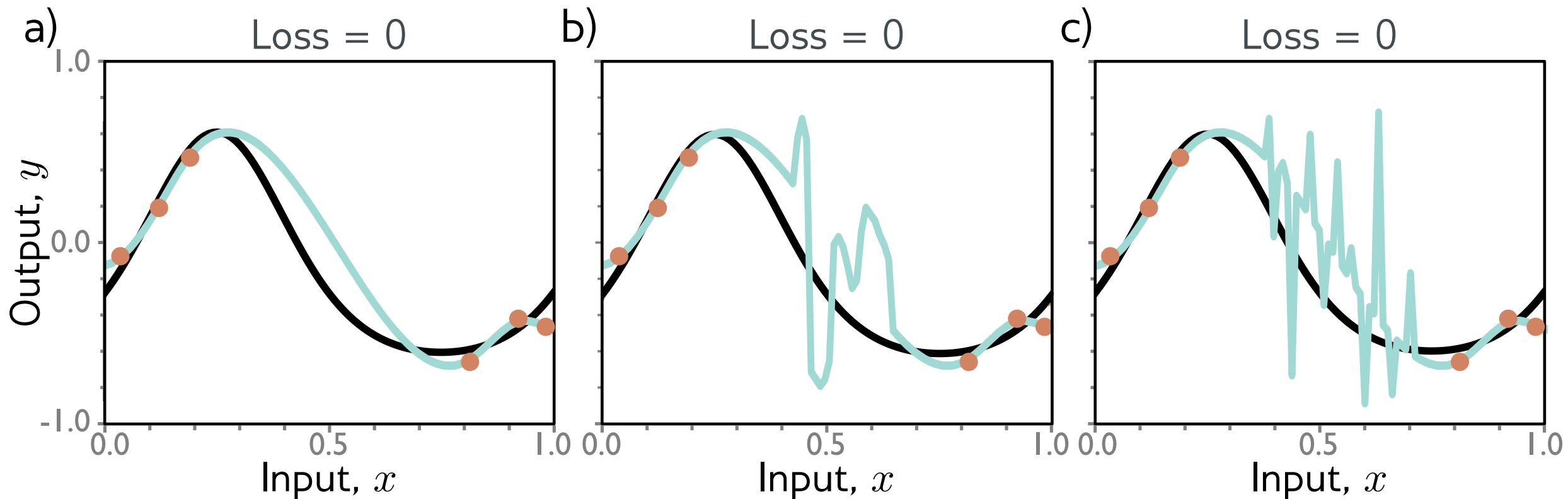
Potential explanation:

- can make smoother functions with more hidden units
- being smooth between the datapoints is a reasonable thing to do

But why?



- All of these solutions are equivalent in terms of loss.
- Why should the model choose the smooth solution?
- Tendency of model to choose one solution over another is **inductive bias**



Possible Explanations:

- Network initialization (He init) encourages smoothness (prior over weights)
 - sensible magnitudes that create smooth functions
- Training procedure is regularized enough to encourage smoothness
 - SGD, weight decay, batch- and layer-norm, dropout, etc...

Measuring performance

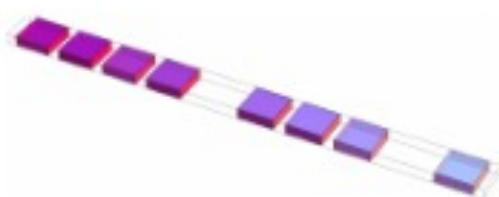
- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Curse of dimensionality & weird properties of high dimensional space
- Choosing hyperparameters

Curse of dimensionality

- 40-dimensional data
- 10,000 data points
- Consider quantizing each dimension into 10 bins
- 10^{40} bins
- 1 data point per 10^{35} bins
- The tendency of high-dimensional space to overwhelm the number of data points is called the **curse of dimensionality**

Curse of dimensionality

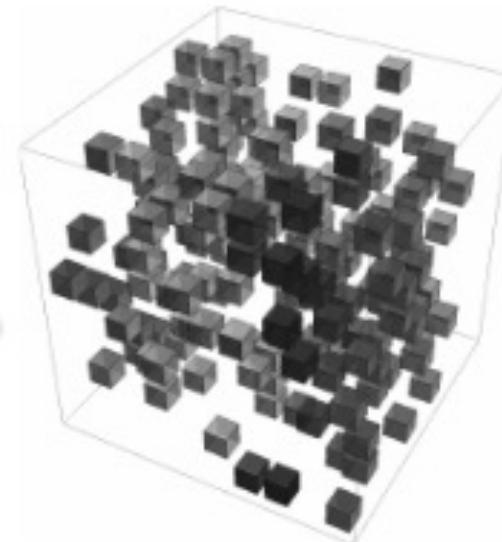
$$10 \times 1 = 10$$



$$10 \times 10 = 100$$



$$10 \times 10 \times 10 = 1000$$



As the number of dimensions grows, the number of possible “bins” grows exponentially. In order to maintain a fixed amount of coverage of bins (here 80%), the amount of data needed also grows exponentially

Measuring performance

- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Curse of dimensionality & weird properties of high dimensional space
- Choosing hyperparameters

Choosing hyperparameters

- Don't know bias or variance
- Don't know how much capacity to add
- How do we choose capacity in practice?
 - Or model structure
 - Or training algorithm
 - Or learning rate
- Third data set – validation set
 - Train models with different hyperparameters on training set
 - Choose best hyperparameters with validation set
 - Test once with test set

Assessing Generalization: Beyond the IID Assumption

IID Recap

Independent

In most of supervised learning problems, we make the assumption that our data are independent and identically distributed (IID)

This means that all samples are independent of each other, and follow the same distribution

This applies to both the **training data** and the **test data**

Gives us confidence that the outputs of our network will generalize



$$Y_i | X_i \perp\!\!\!\perp Y_j | x_j$$

IID Recap

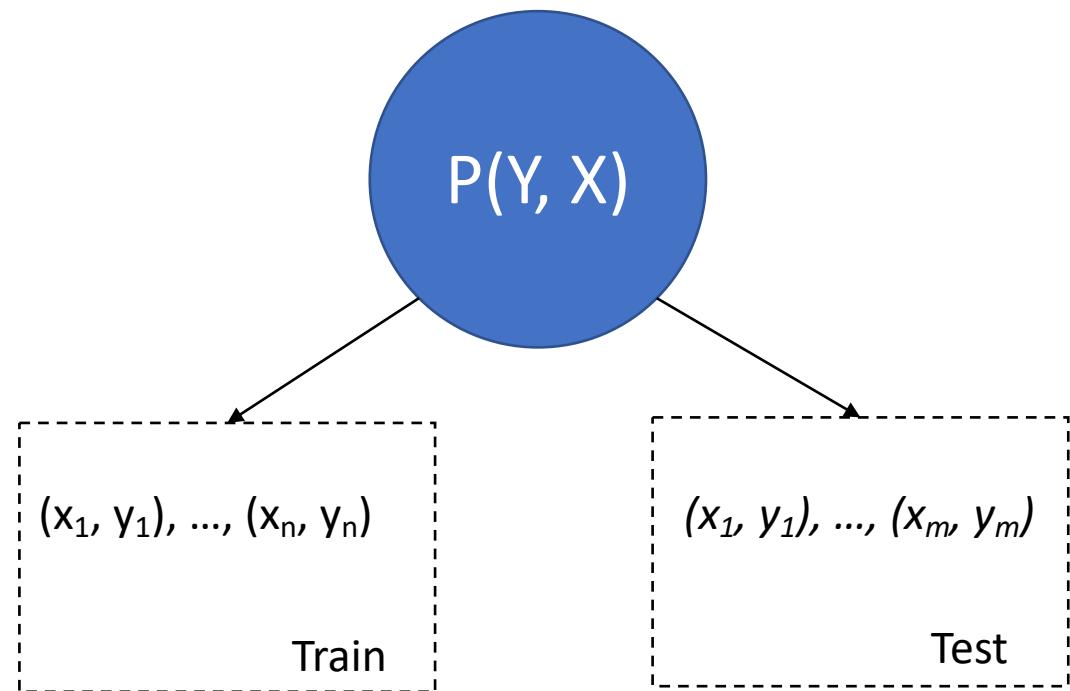
Identically distributed

In most of supervised learning problems, we make the assumption that our data are independent and identically distributed (IID)

This means that all samples are independent of each other, and follow the same distribution

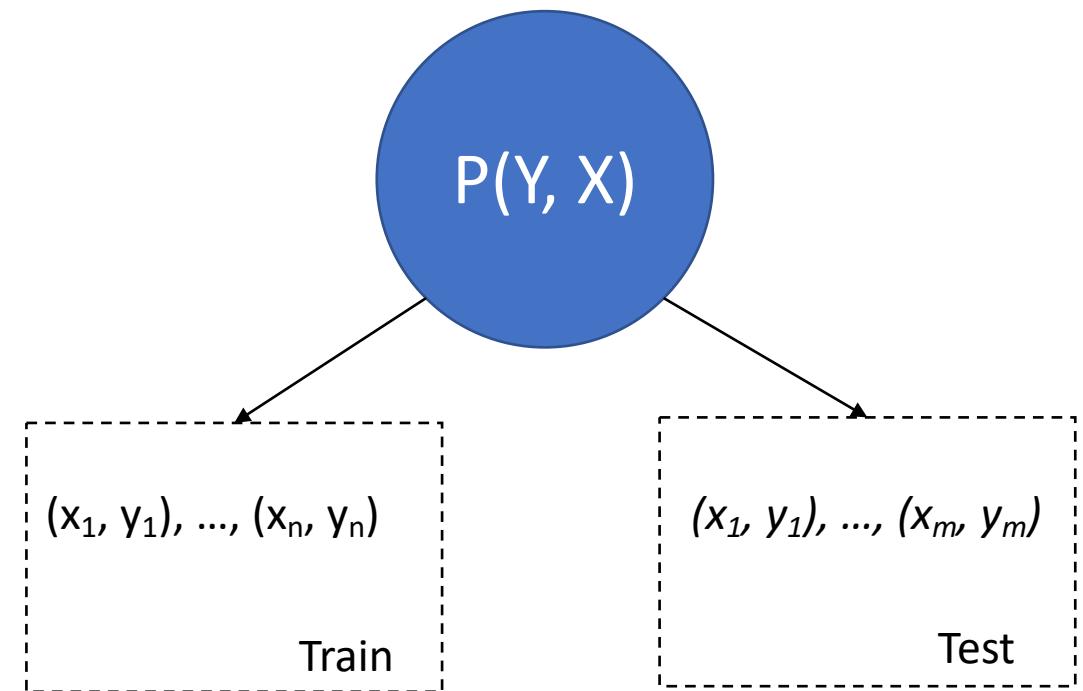
This applies to both the **training data** and the **test data**

Gives us confidence that the outputs of our network will generalize



IID Recap

Identically distributed



Allows us to learn model
parameters via maximum
likelihood

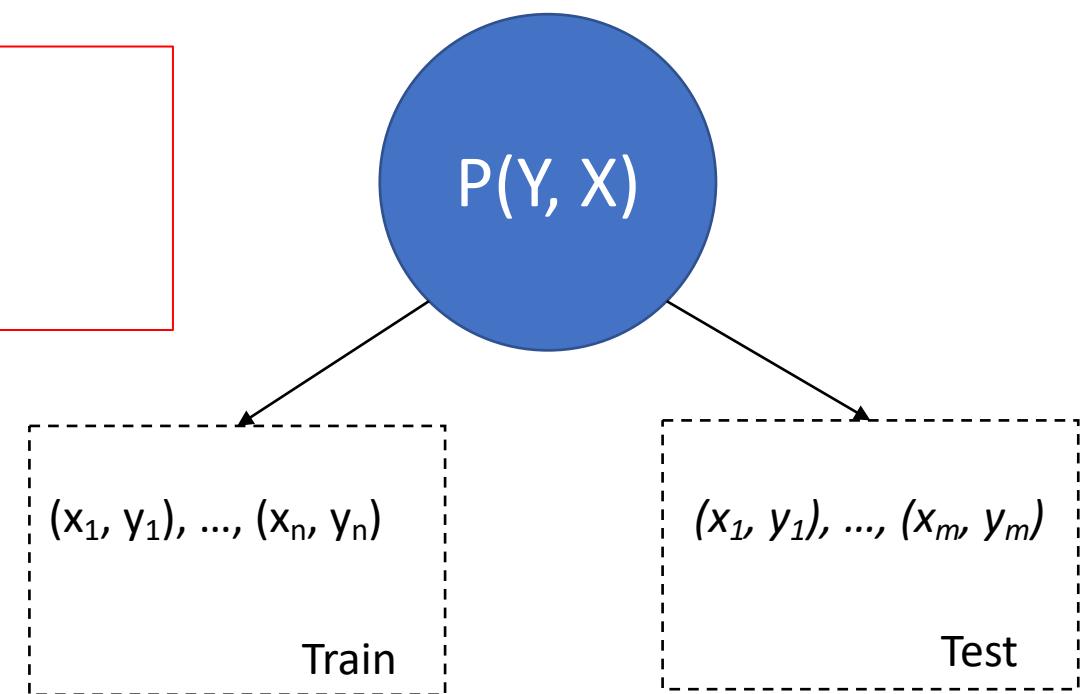
IID Recap

Identically distributed

$$\hat{\phi} = \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right]$$

Product Rule of Probability:

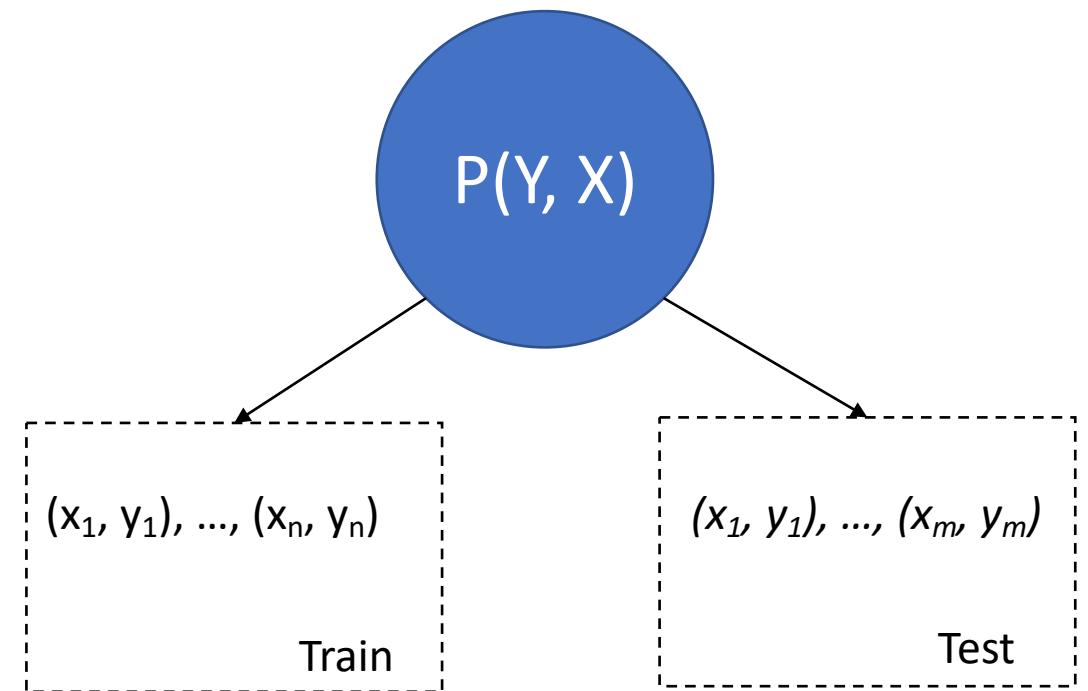
If A and B are independent:
then $P(A \text{ and } B) = P(A) * P(B)$



Allows us to learn model
parameters via maximum
likelihood

IID Recap

Identically distributed



Allows us to generalize to
new data

Distribution Shift

Unfortunately, modern machine learning is not so clean...

We often observe shifts of the distribution between when we train, and when the model is actually performing inference

Neural networks are actually quite bad at adapting to distributional changes

Distribution Shift

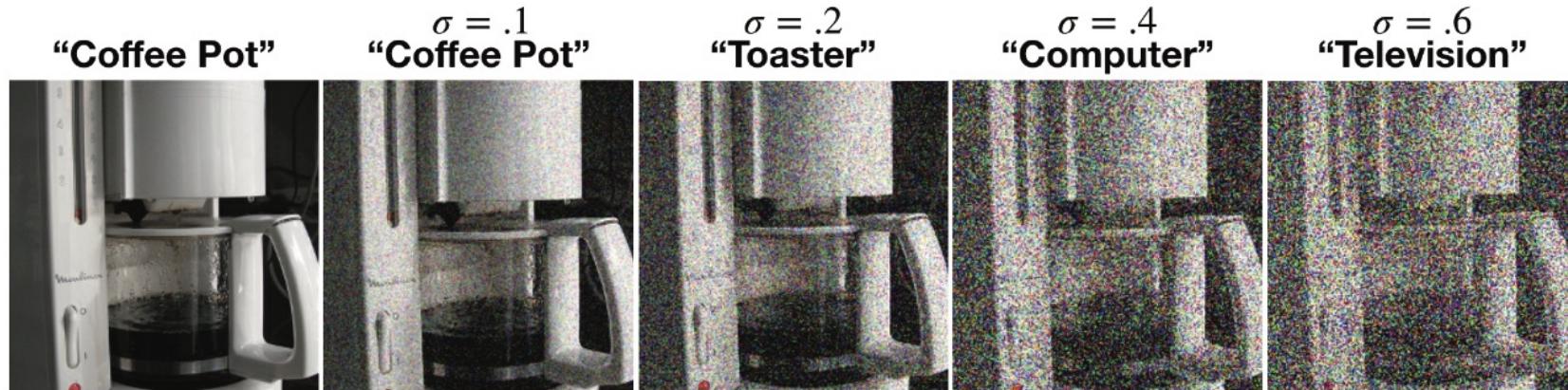


Figure 19.1: Effect of Gaussian noise of increasing magnitude on an image classifier. The model is a ResNet-50 CNN trained on ImageNet. From Figure 23 of [For+19]. Used with kind permission of Justin Gilmer.

Distribution Shift



(A) **Cow: 0.99**, Pasture: 0.99, Grass: 0.99, No Person: 0.98, Mammal: 0.98



(B) No Person: 0.99, Water: 0.98, Beach: 0.97, Outdoors: 0.97, Seashore: 0.97



(C) No Person: 0.97, **Mammal: 0.96**, Water: 0.94, Beach: 0.94, Two: 0.94

Distribution Shift



(A) **Cow: 0.99**, Pasture: 0.99, Grass: 0.99, No Person: 0.98, Mammal: 0.98



(B) No Person: 0.99, Water: 0.98, Beach: 0.97, Outdoors: 0.97, Seashore: 0.97



(C) No Person: 0.97, **Mammal: 0.96**, Water: 0.94, Beach: 0.94, Two: 0.94

Presence of grass is a spurious correlation: predictive of examples in the test set, doesn't generalize

Distribution Shift

High dimensional analogue of this:



(A) **Cow: 0.99**, Pasture: 0.99, Grass: 0.99, No Person: 0.98, Mammal: 0.98



(B) No Person: 0.99, Water: 0.98, Beach: 0.97, Outdoors: 0.97, Seashore: 0.97



(C) No Person: 0.97, **Mammal: 0.96**, Water: 0.94, Beach: 0.94, Two: 0.94



Presence of grass is a spurious correlation: predictive of examples in the test set, doesn't generalize

Distribution Shift – Motivating Examples

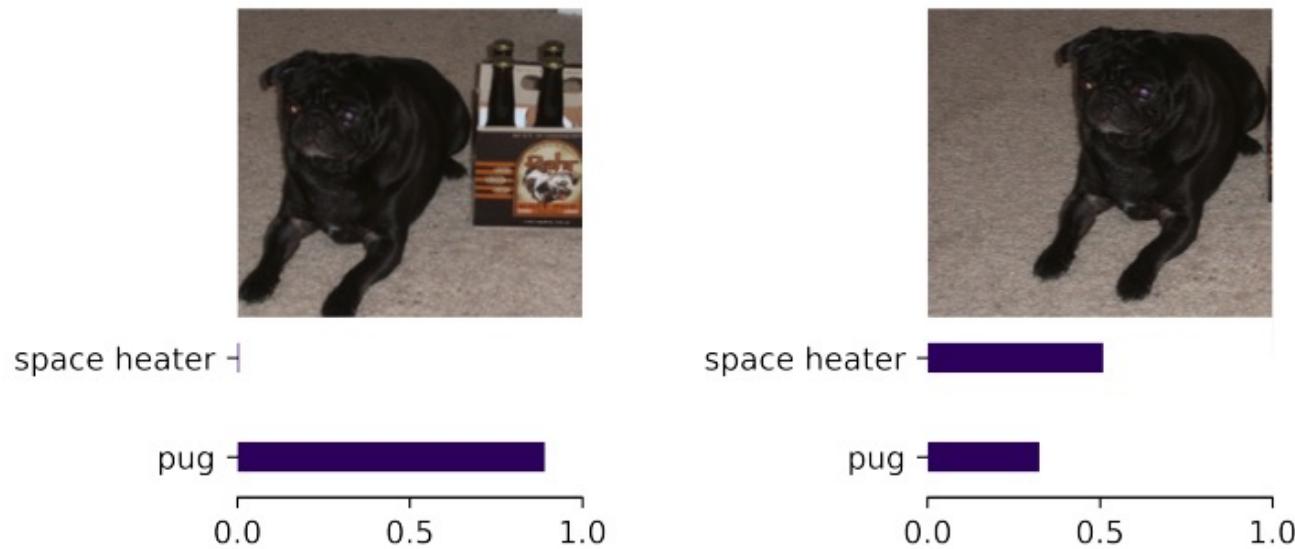


Figure 5. A correctly and an incorrectly classified sample for the BEiT-Model. The dog of class “pug” is fully visible in both images.

Causal Definition of Distribution Shift

Inputs X and outputs Y

Definition 1: Discriminative Model

Data generating process: $X \rightarrow Y$

Example: Suppose X is a medical image and Y is segmentation created by a doctor. The image is causal of the segmentation. If we change X , then Y would also change Y

$$p_{\theta}(x, y) = p_{\psi}(x)p_w(y|x)$$

Causal Definition of Distribution Shift

Inputs X and outputs Y

Definition 2: Generative Model

Data generating process: $Y \rightarrow X$

Example: Suppose X is a medical image and Y is the ground truth disease state

Changing the disease state, Y , will change the appearance of the image

$$p_{\theta}(x, y) = p_{\pi}(y)p_{\phi}(x|y)$$

Covariate Shift

Discriminative Model: $p_{\theta}(x, y) = p_{\psi}(x)p_w(y|x)$

Consider if we change the distribution of the inputs $p_{\psi}(x)$ from the source distribution (train) to the target distribution (test)

$$\psi^s \neq \psi^t$$

Covariate Shift

Discriminative Model: $p_{\theta}(x, y) = p_{\psi}(x)p_w(y|x)$

Consider if we change the distribution of the inputs $p_{\psi}(x)$ from the source distribution (train) to the target distribution (test)

$$\psi^s \neq \psi^t$$

Example: We train an object detector for products in marketing images and test on products on cluttered shelves

Source:



Target:



Concept Shift

Discriminative Model: $p_{\theta}(x, y) = p_{\psi}(x)p_w(y|x)$

Consider if we change the distribution of the conditional relationship between x and y , $p_w(y|x)$ from the source distribution (train) to the target distribution (test)

$$w^s \neq w^t$$

Example: Different doctors use different rules and heuristics to label pixels of a medical image

Label Shift

Generative Model: $p_{\theta}(x, y) = p_{\pi}(y)p_{\phi}(x|y)$

Consider if we change the distribution of the outputs, $p_{\pi}(y)$ from the source distribution (train) to the target distribution (test)

$$\phi^s \neq \phi^t$$

Example: Suppose that the medical images about a disease are collected in an urban source environment, but the target environment is in rural areas

The prevalence of the disease may be different between the two

Manifestation Shift

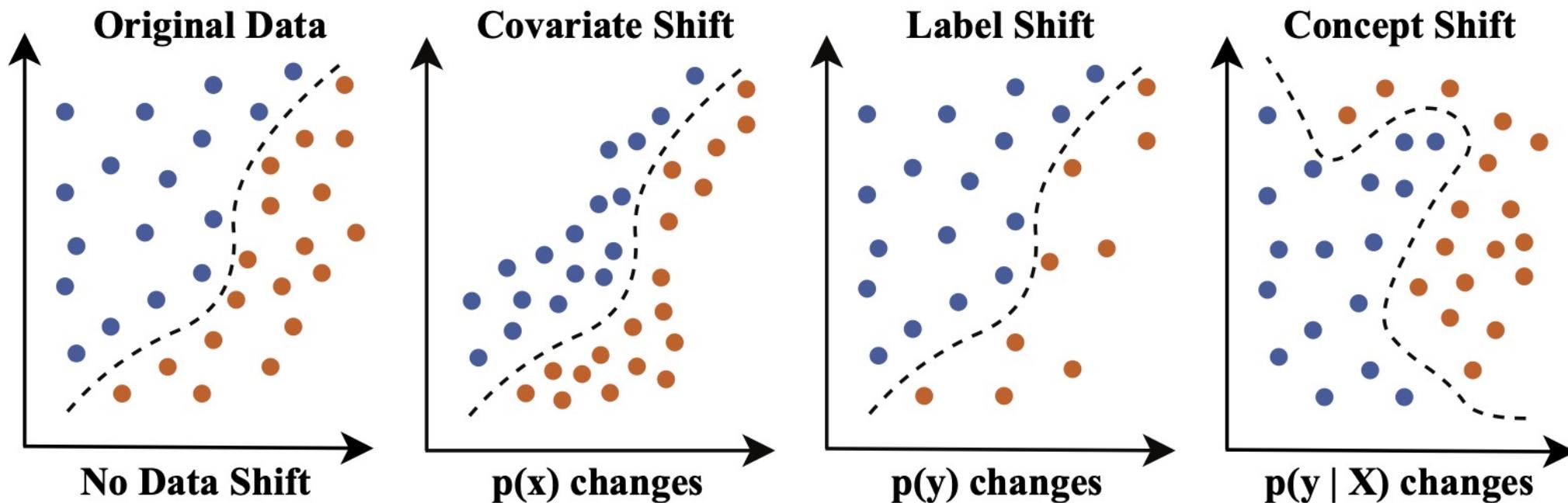
Generative Model: $p_{\theta}(\mathbf{x}, \mathbf{y}) = p_{\pi}(\mathbf{y})p_{\phi}(\mathbf{x}|\mathbf{y})$

Consider if we change the distribution of the inverse mapping, $p_{\phi}(\mathbf{x}|\mathbf{y})$ from the source distribution (train) to the target distribution (test)

$$\phi^s \neq \phi^t$$

Example: The way that a disease, Y, manifests itself in image, X, changes maybe due to some other hidden variable, Z, such as age.

Different Types of Distribution Shift



Different Types of Distribution Shift

Name	Source	Target	Joint
Covariate/domain shift	$p(X)p(Y X)$	$q(X)p(Y X)$	Discriminative
Concept shift	$p(X)p(Y X)$	$p(X)q(Y X)$	Discriminative
Label (prior) shift	$p(Y)p(X Y)$	$q(Y)p(X Y)$	Generative
Manifestation shift	$p(Y)p(X Y)$	$p(Y)q(X Y)$	Generative

Table 19.1: The 4 main types of distribution shift.

Selection Bias

How we collect the data (train and test) can impact how our model performances

- Suppose we want train an image classifier for a production software system.
- During training we filter to clean, non-blurry images
- At test time, we don't have this constraint

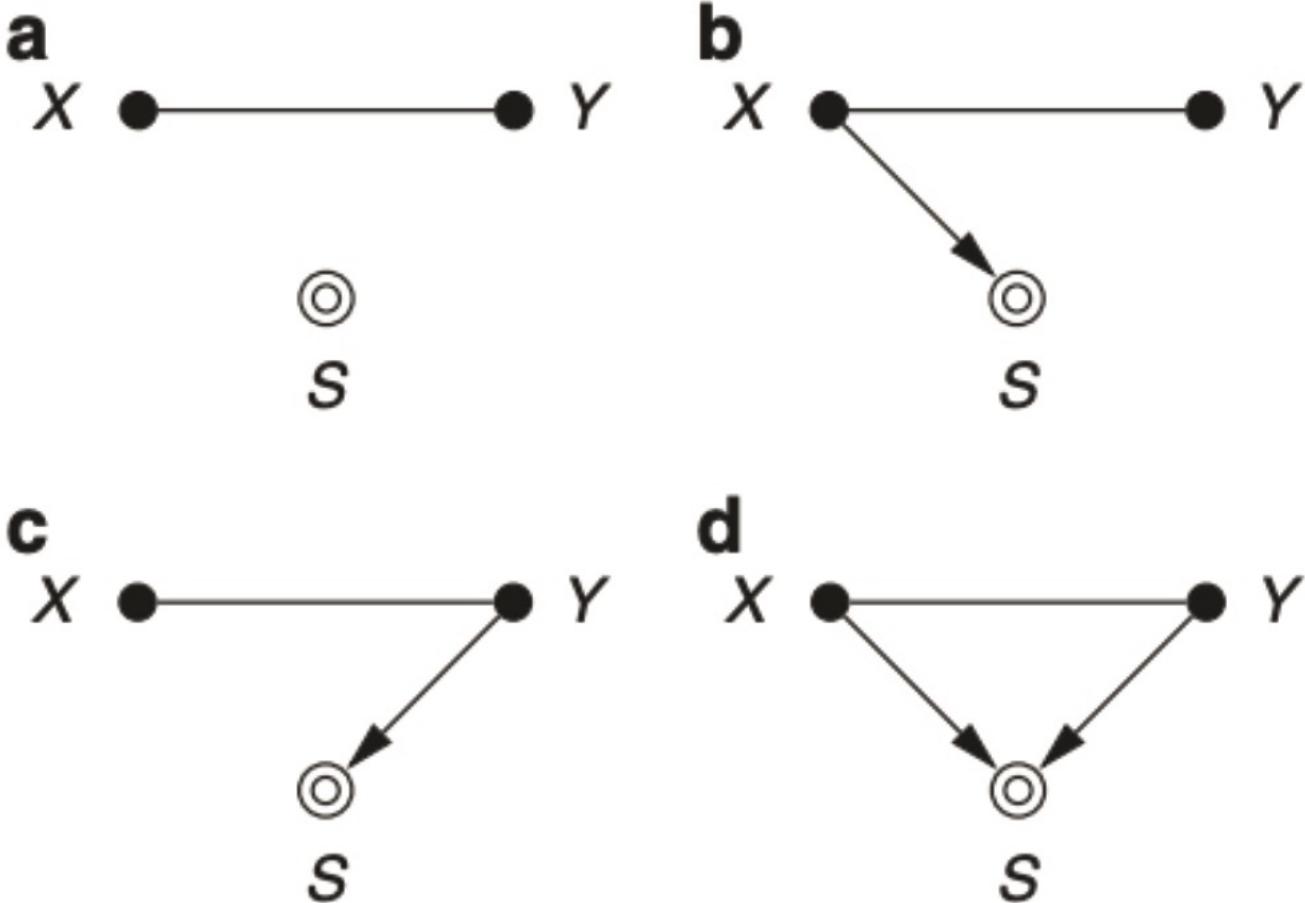
The selection mechanism for the data induces **covariate shift**

Selection Bias

Input: X

Label: Y,

S: Selection Variable {0, 1}



Selection Bias

Input: X

Label: Y,

S: Selection Variable {0, 1}

No selection

a



S

b



S

c



S

d



S

Selection Bias

Input: X

Label: Y,

S: Selection Variable {0, 1}

a



S

Selection
based on X

b



S

c



S

d



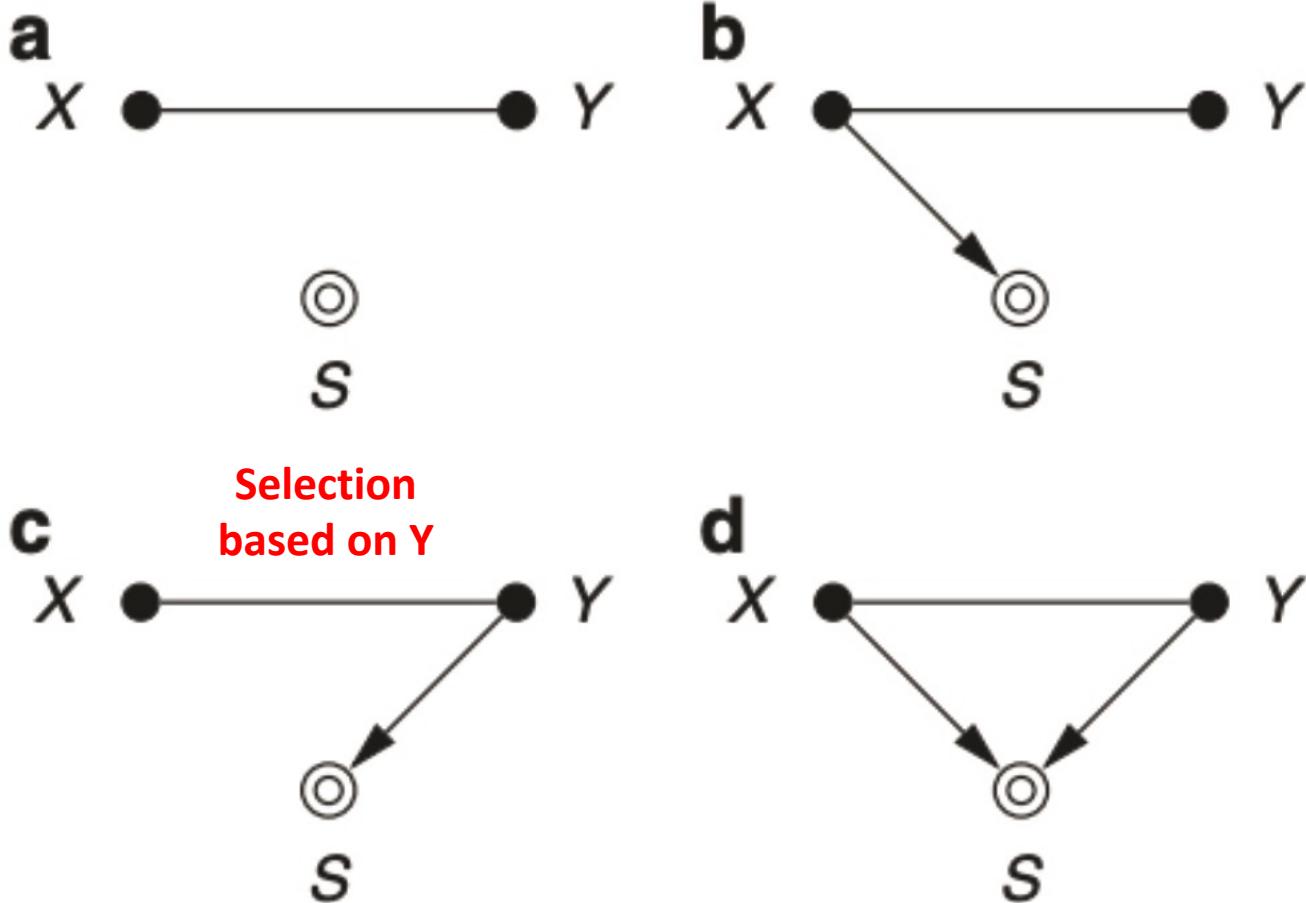
S

Selection Bias

Input: X

Label: Y,

S: Selection Variable {0, 1}

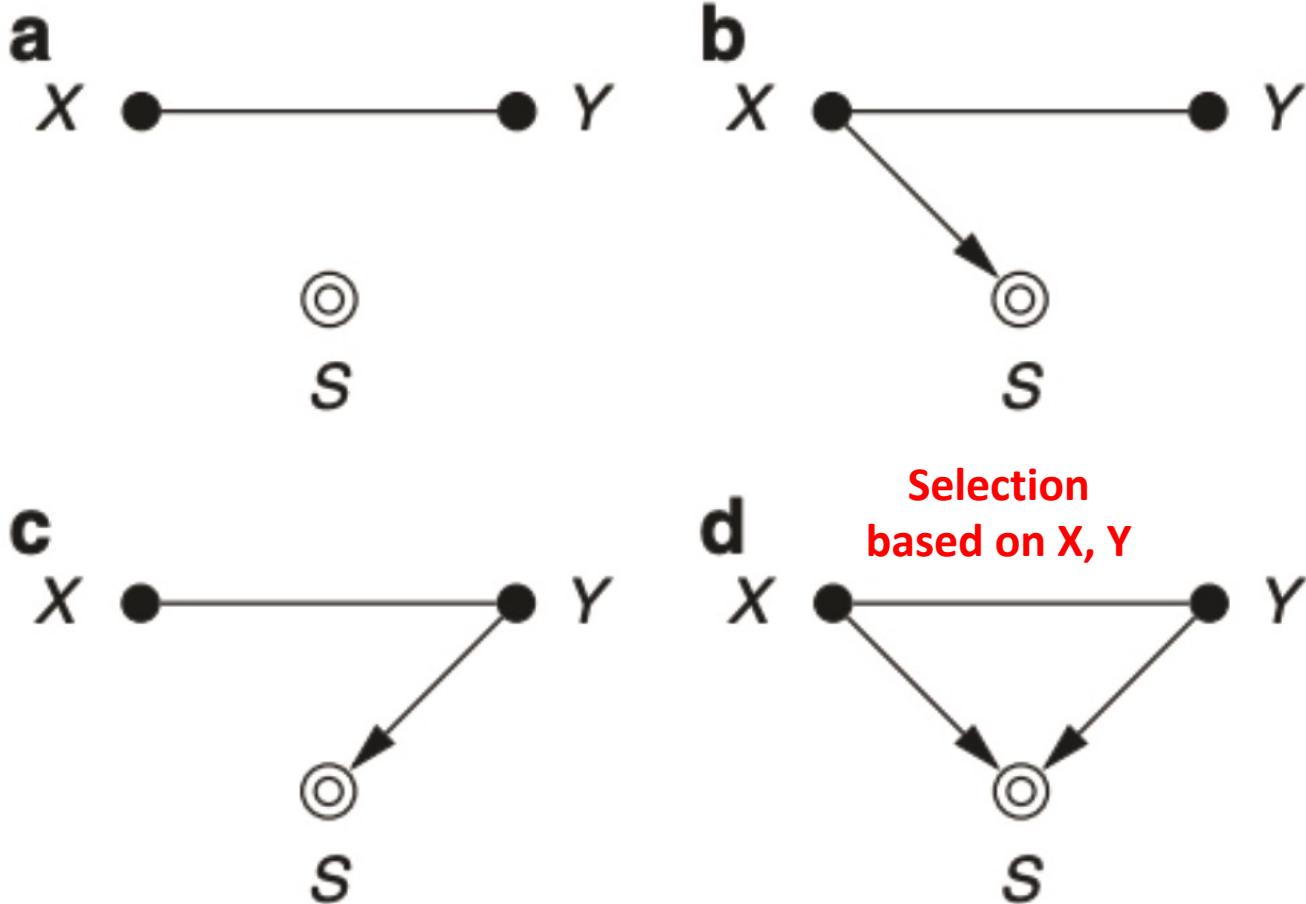


Selection Bias

Input: X

Label: Y,

S: Selection Variable {0, 1}



Detecting Shifts

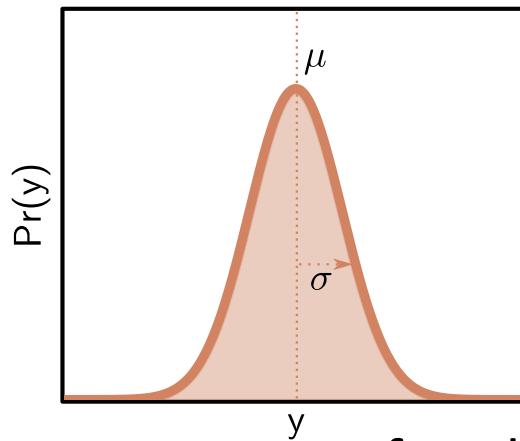
- A variety of ideas exist to detect out-of-distribution (OOD).
- One good idea is to use networks that output a probability distribution, rather than models that output a single point prediction

How can a model predict a probability distribution?

1. Pick a known distribution (e.g., normal distribution) to model output y with parameters θ

e.g., the normal distribution

$$\theta = \{\mu, \sigma^2\}$$



2. Use model to predict parameters θ of probability distribution
 - $y_i = f(x_i; \varphi)$
 - $\mu_i, \sigma_i = f(x_i; \varphi)$

Heteroscedastic regression

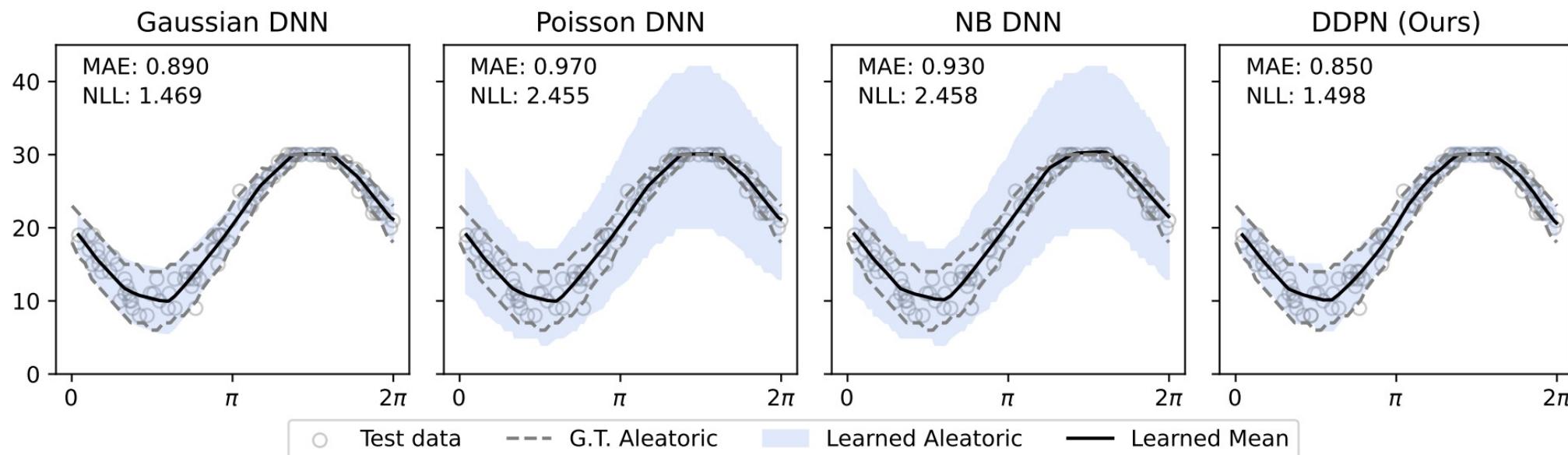
- Previously assumed that the noise σ^2 is the same everywhere.
- But we could make the noise a function of the data x .
- Build a model with two outputs:

$$\mu = f_1[\mathbf{x}, \phi]$$

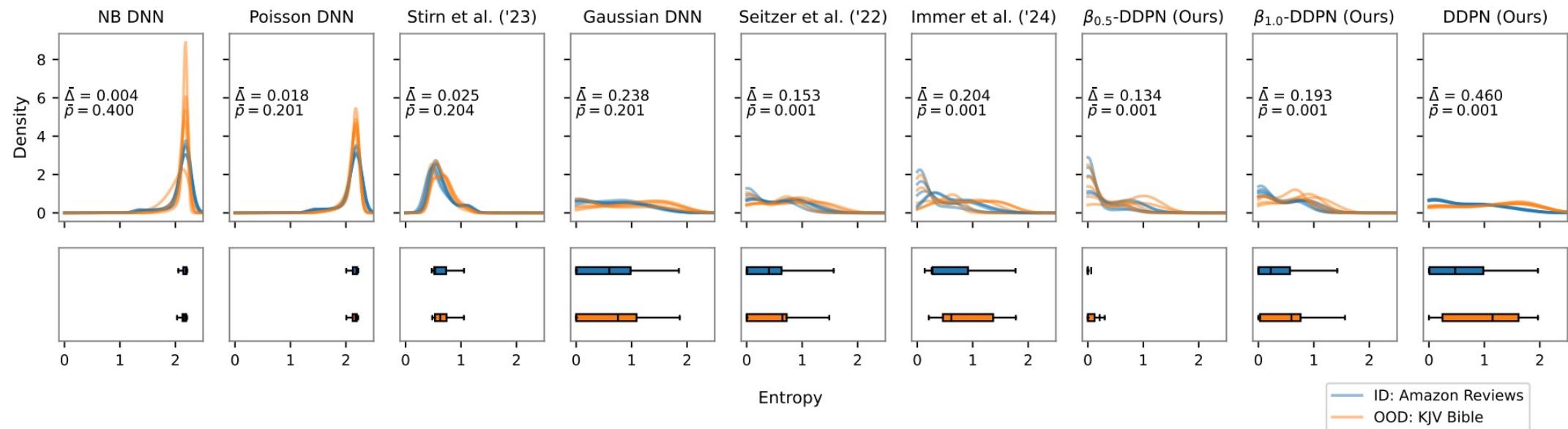
$$\sigma^2 = f_2[\mathbf{x}, \phi]^2$$

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I \log \left[\frac{1}{\sqrt{2\pi f_2[\mathbf{x}_i, \phi]^2}} \right] - \frac{(y_i - f_1[\mathbf{x}_i, \phi])^2}{2f_2[\mathbf{x}_i, \phi]^2} \right]$$

Detecting Shifts



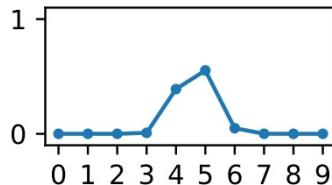
Entropy of Probabilistic Models



Entropy of Probabilistic Models

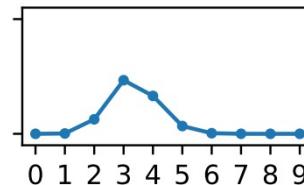
DDPN (Ours)

"Plant was small but in good shape. After repotting right off the bat, and 4 months of intermittent neglect, it's still going strong."



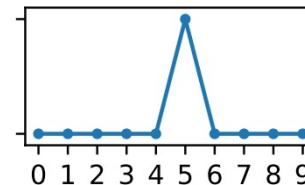
"And to Seth, to him also there was born a son; and he called his name Enos: then began men to call upon the name of the LORD."

"Small pump, much smaller than I thought I was going to get. Worked ok for a while but it has to be in a lot of sun."



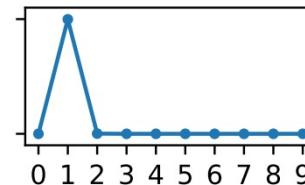
"And he left off talking with him, and God went up from Abraham."

"Best gas can I ever owned. No-spill, or smell in the service van."

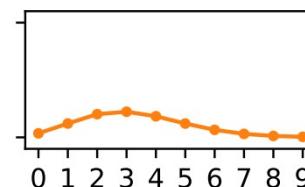
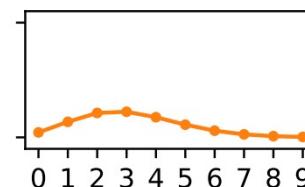
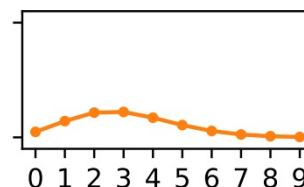
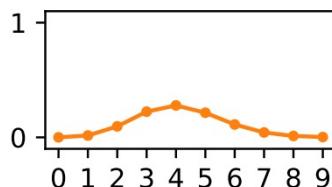


"And Moses returned unto the LORD, and said, Lord, wherefore hast thou so evil entreated this people? why is it that thou hast sent me?"

"Bought this item Dec 2017, were now at April 2018, so 4 months and now it sinks to the bottom of the pool, will not float anymore and my return window closed Feb 2018. Stay away, what a waste of money!!!"



"And they pitched by Jordan, from Bethjesimoth even unto Abel-shittim in the plains of Moab."



Questions?