

# Deep Learning for Data Science

## DS 542

<https://dl4ds.github.io/fa2025/>

Measuring Performance and Generalization



# Plan for Today

- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Choosing hyperparameters

# MNIST1D

---

## Scaling down Deep Learning

---

Sam Greydanus<sup>1</sup>

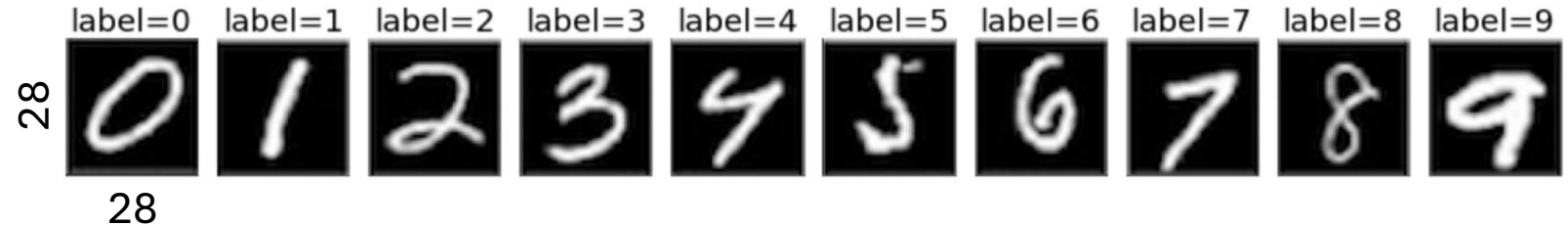
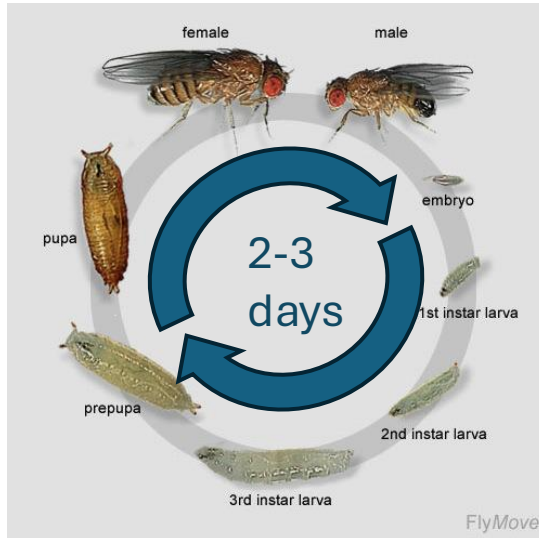
“A large number of deep learning innovations including [dropout](#), [Adam](#), [convolutional networks](#), [generative adversarial networks](#), and [variational autoencoders](#) began life as MNIST experiments. Once these innovations proved themselves on small-scale experiments, scientists found ways to scale them to larger and more impactful applications.”

S. Greydanus, “Scaling down Deep Learning.” arXiv, Dec. 04, 2020. doi: [10.48550/arXiv.2011.14439](https://arxiv.org/abs/10.48550/arXiv.2011.14439).

<https://github.com/greydanus/mnist1d>

# MNIST Dataset

- 28x28x1 grayscale images
- 60K Training, 10K Test
- “Is to Deep Learning what fruit flies are to genetics research”

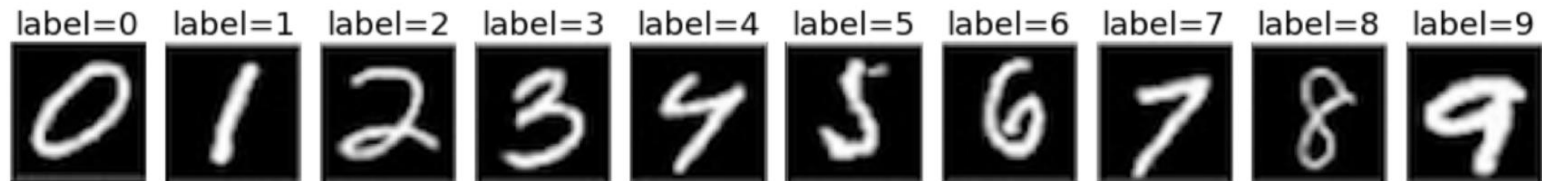


But poorly differentiates model performance:

Model Type	Accuracy
Logistic Regression	94%
MLP	99+%
CNN	99+%

# MNIST 1D Dataset

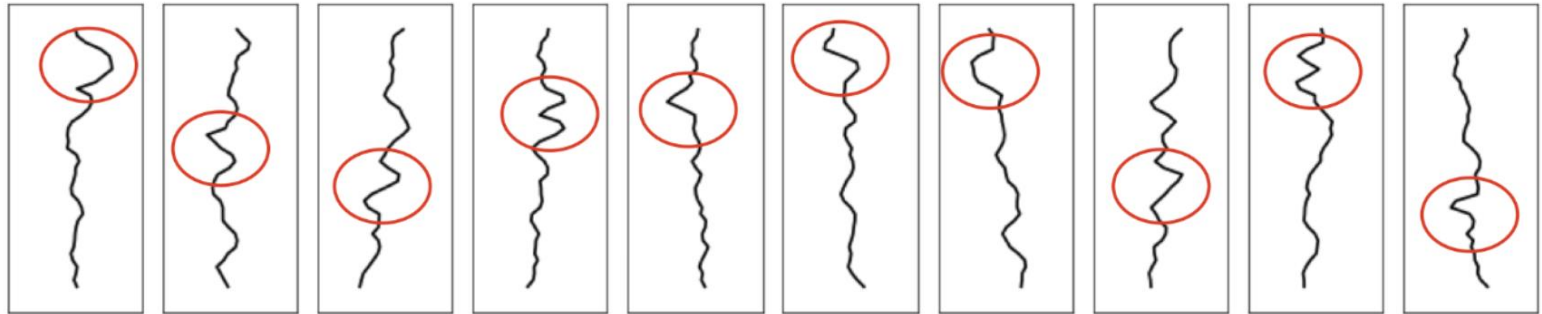
Original MNIST examples



Represent digits as 1D patterns



Pad, translate & transform



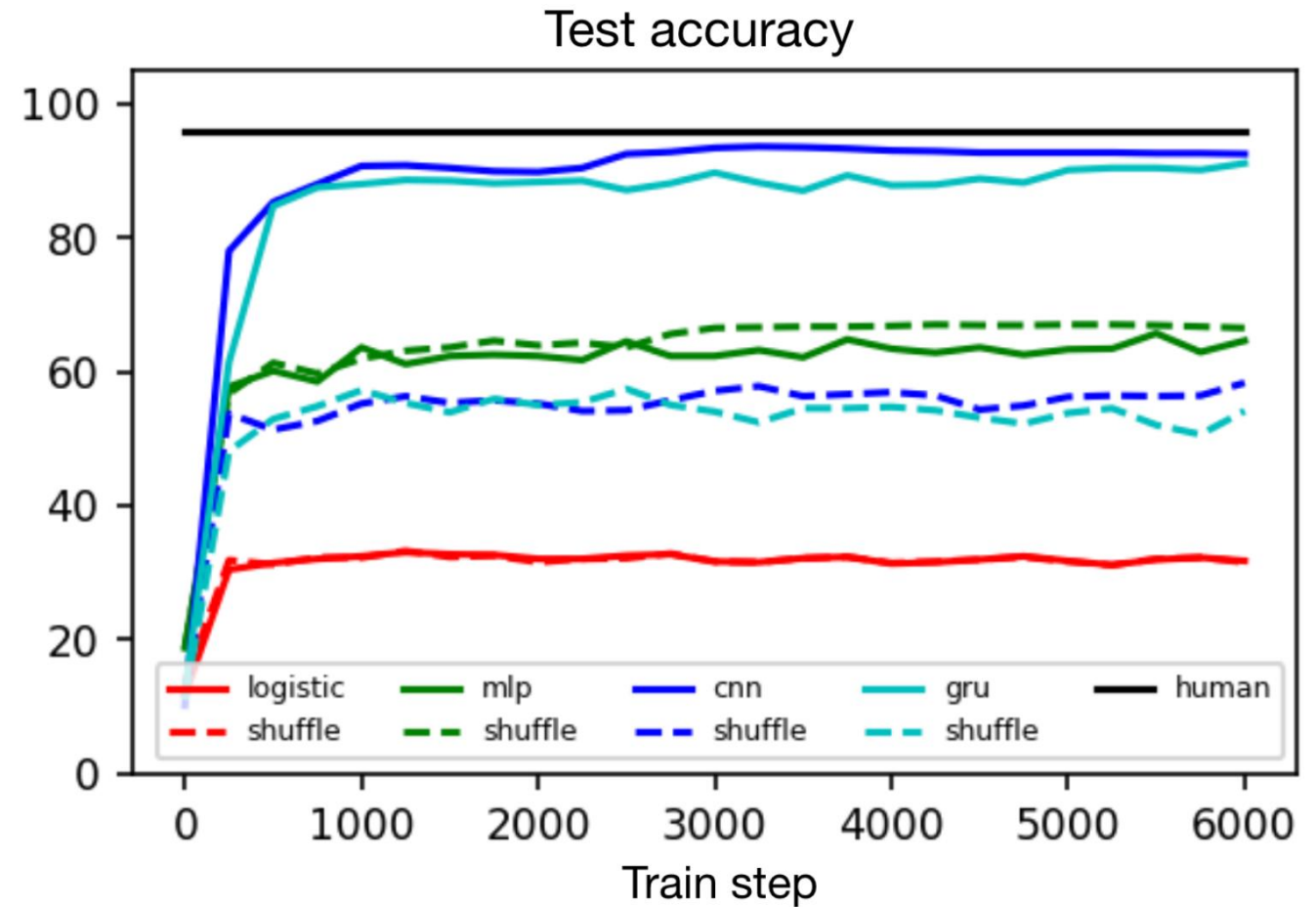
Fixed, 1-D, length-12  
templates for each of 10  
digit classes

Generate dataset by  
programmatically applying  
6 parametric  
transformations.

E.g. pad, shear, translate, correlated noise, i.i.d. noise, interpolation.

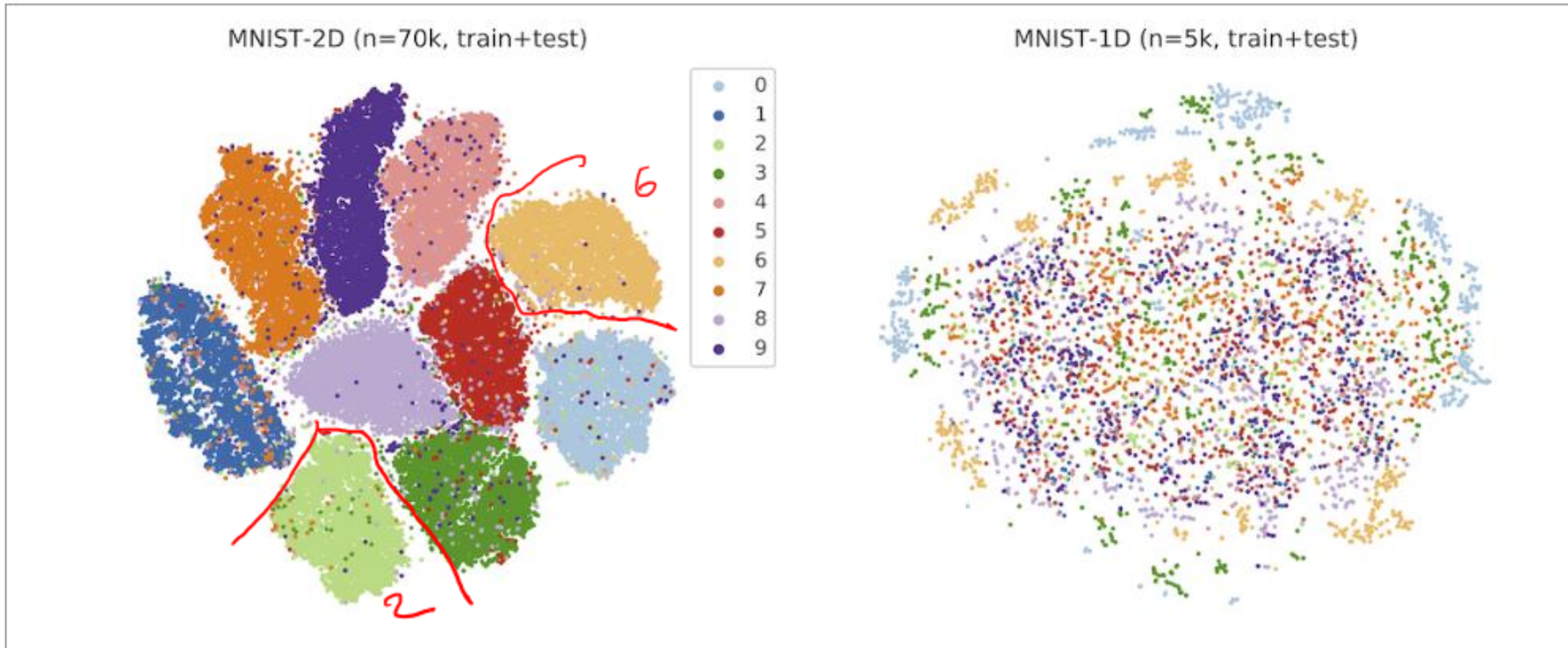
# MNIST 1D

Differentiates performance of different model types much more than MNIST



Dataset	Logistic regression	Fully connected model	Convolutional model	GRU model	Human expert
MNIST	$94 \pm 0.5$	$> 99$	$> 99$	$> 99$	$> 99$
MNIST-1D	$32 \pm 1$	$68 \pm 2$	$94 \pm 2$	$91 \pm 2$	$96 \pm 1$
MNIST-1D (shuffled)	$32 \pm 1$	$68 \pm 2$	$56 \pm 2$	$57 \pm 2$	$\approx 30 \pm 10$

# Visualizing MNIST and MNIST-1D with tSNE



Visualizing the MNIST and MNIST-1D datasets with tSNE. The well-defined clusters in the MNIST plot indicate that the majority of the examples are separable via a kNN classifier in pixel space. The MNIST-1D plot, meanwhile, reveals a lack of well-defined clusters which suggests that learning a nonlinear representation of the data is much more important to achieve successful classification. Thanks to [Dmitry Kobak](#) for making this plot.

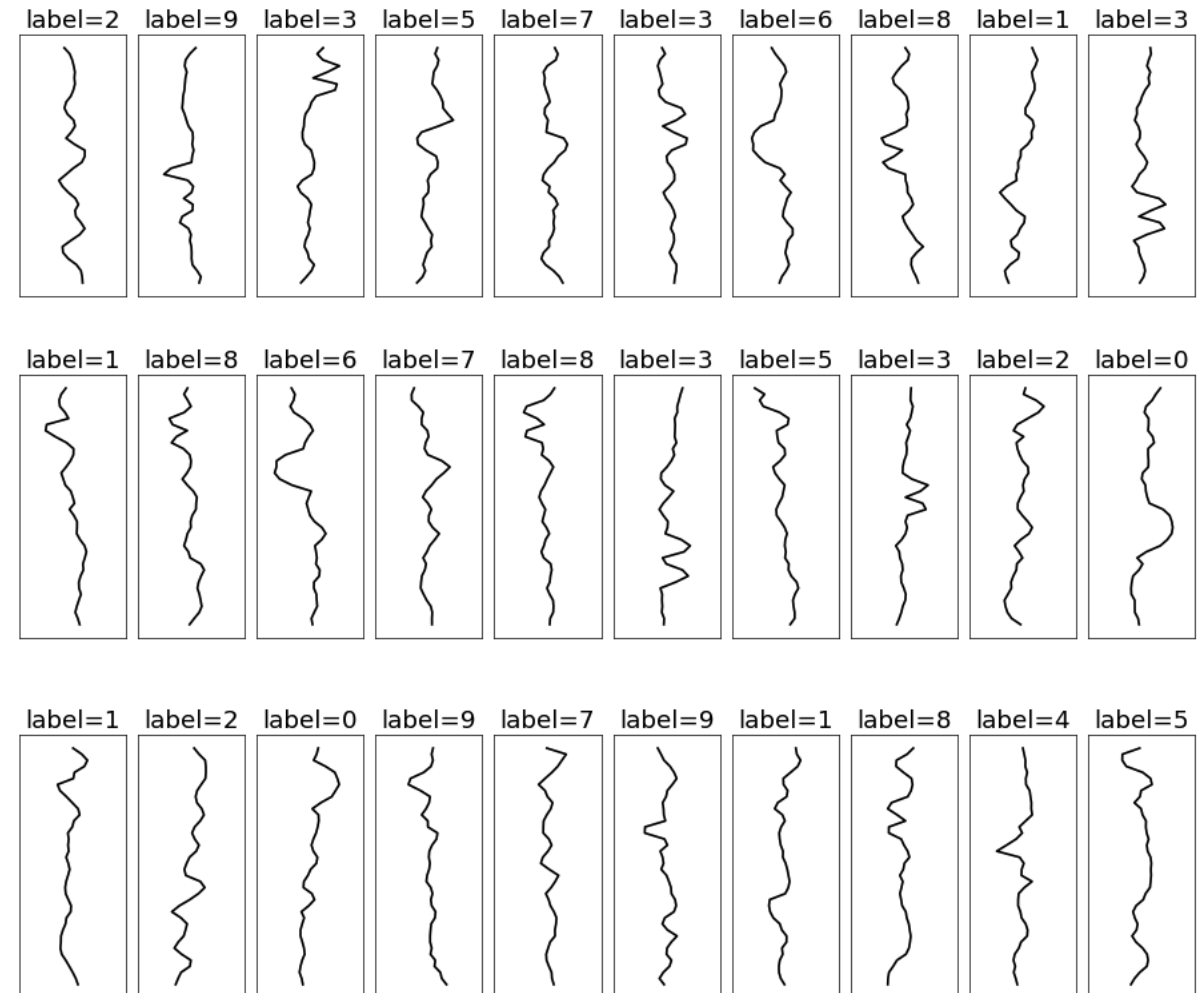
<https://twitter.com/hippopedoid>



# MNIST1D Train and Test Set

- 1D, Length 40 samples
- 4,000 training samples
- 1,000 test samples (80/20 split)

## Dataset Samples





# Network

→ 40 inputs

- 10 outputs

- Two hidden layers
  - 100 hidden units each

- SGD with batch size 100, learning rate 0.1

- 6000 steps (?? Epochs)

```
model = torch.nn.Sequential(  
    torch.nn.Linear(40, 100),  
    torch.nn.ReLU(),  
    torch.nn.Linear(100, 100),  
    torch.nn.ReLU(),  
    torch.nn.Linear(100, 10))
```

```
# choose cross entropy loss function  
loss_function = torch.nn.CrossEntropyLoss()  
  
# construct SGD optimizer and initialize learning rate and momentum  
optimizer = torch.optim.SGD(model.parameters(), lr = 0.1)  
  
# object that decreases learning rate by half every 10 epochs  
scheduler = StepLR(optimizer, step_size=10, gamma=0.5)  
  
# load the data into a class that creates the batches  
data_loader = DataLoader(TensorDataset(x_train,y_train), batch_size=100, shuffle=True)
```

• • •

```
# inference – just choose the max  
pred_train = model(x_train)  
pred_test = model(x_test)  
_, predicted_train_class = torch.max(pred_train.data, 1)  
_, predicted_test_class = torch.max(pred_test.data, 1)
```

```
=====
```

Layer (type:depth-idx)	Output Shape	Param #
Sequential	[1, 10]	–
└─Linear: 1-1	[1, 100]	4,100
└─ReLU: 1-2	[1, 100]	–
└─Linear: 1-3	[1, 100]	10,100
└─ReLU: 1-4	[1, 100]	–
└─Linear: 1-5	[1, 10]	1,010

```
=====
```

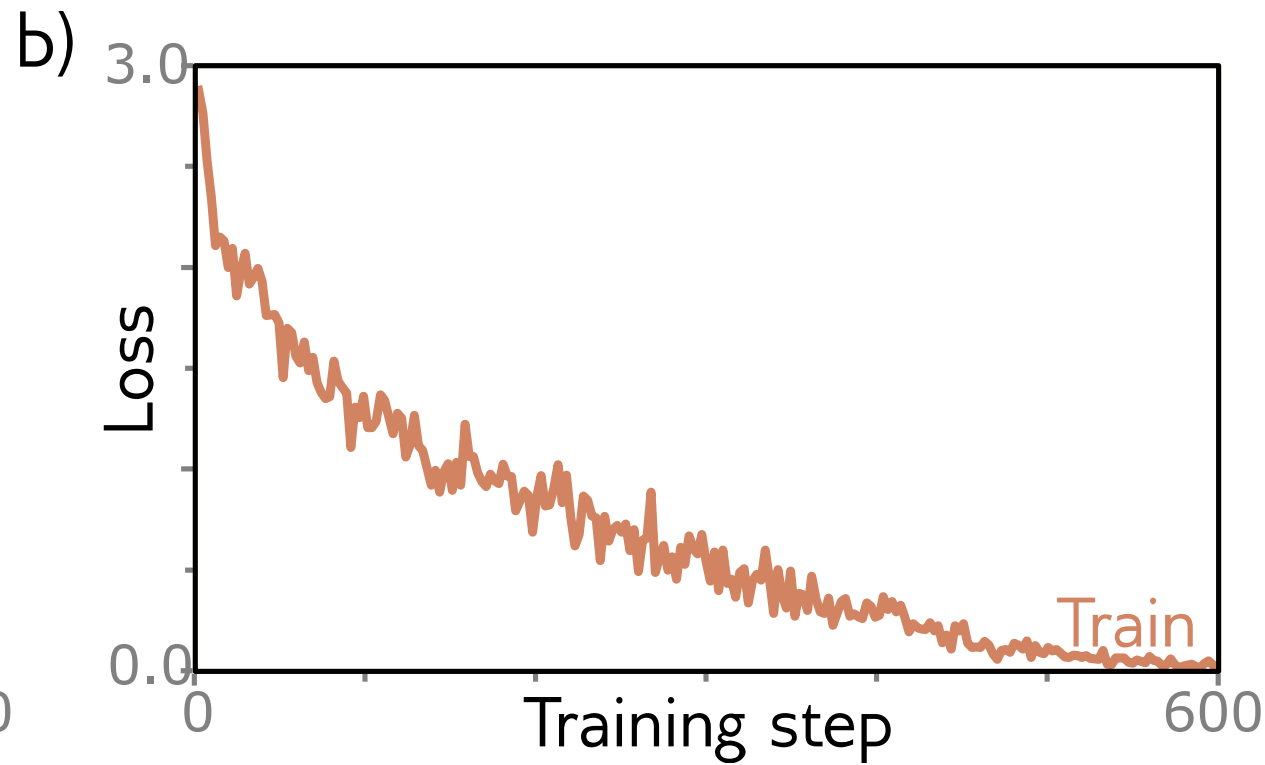
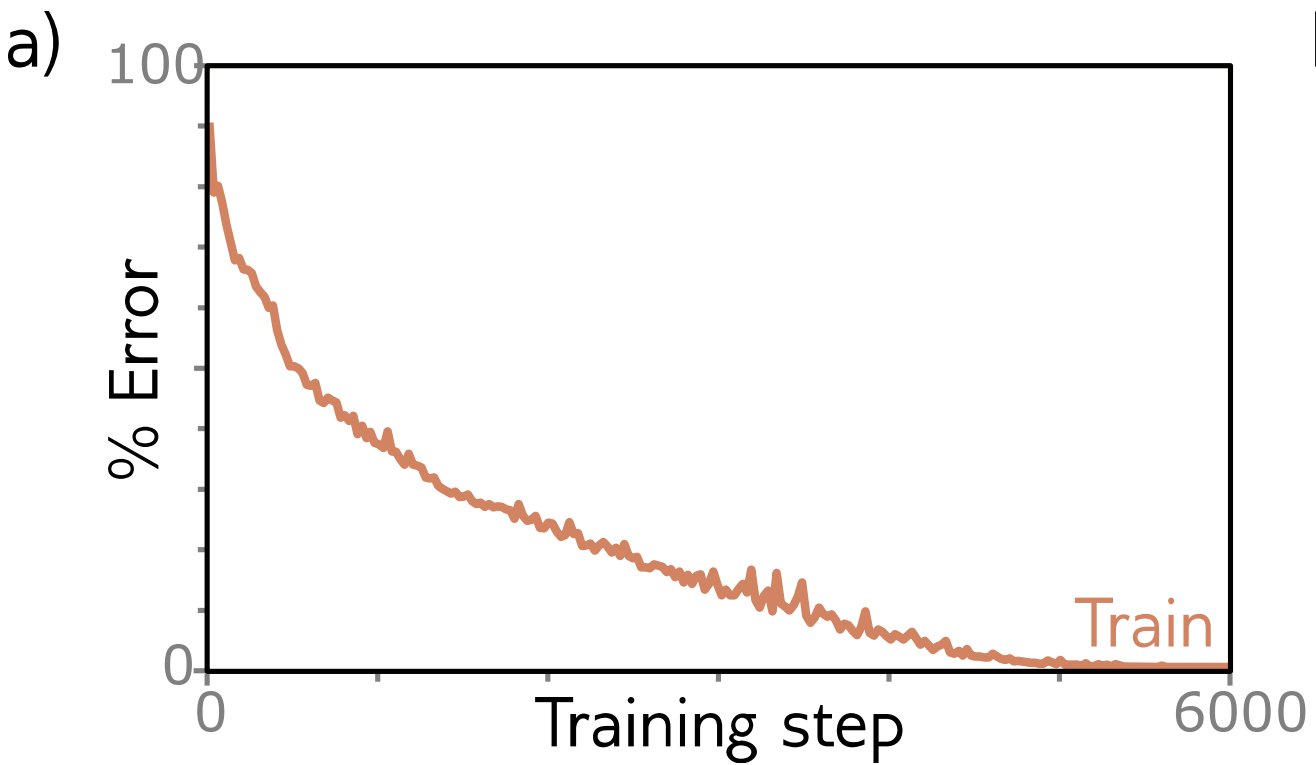
Total params: 15,210  
Trainable params: 15,210  
Non-trainable params: 0  
Total mult-adds (Units.MEGABYTES): 0.02

```
=====
```

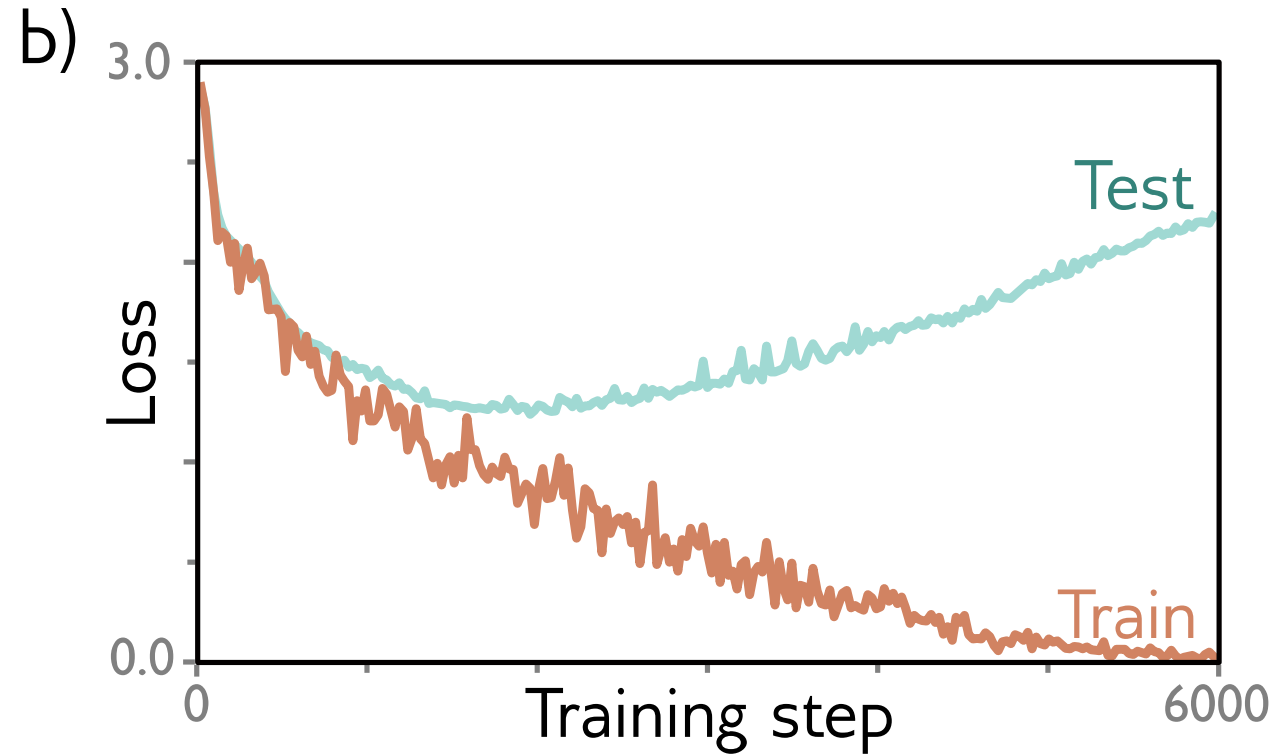
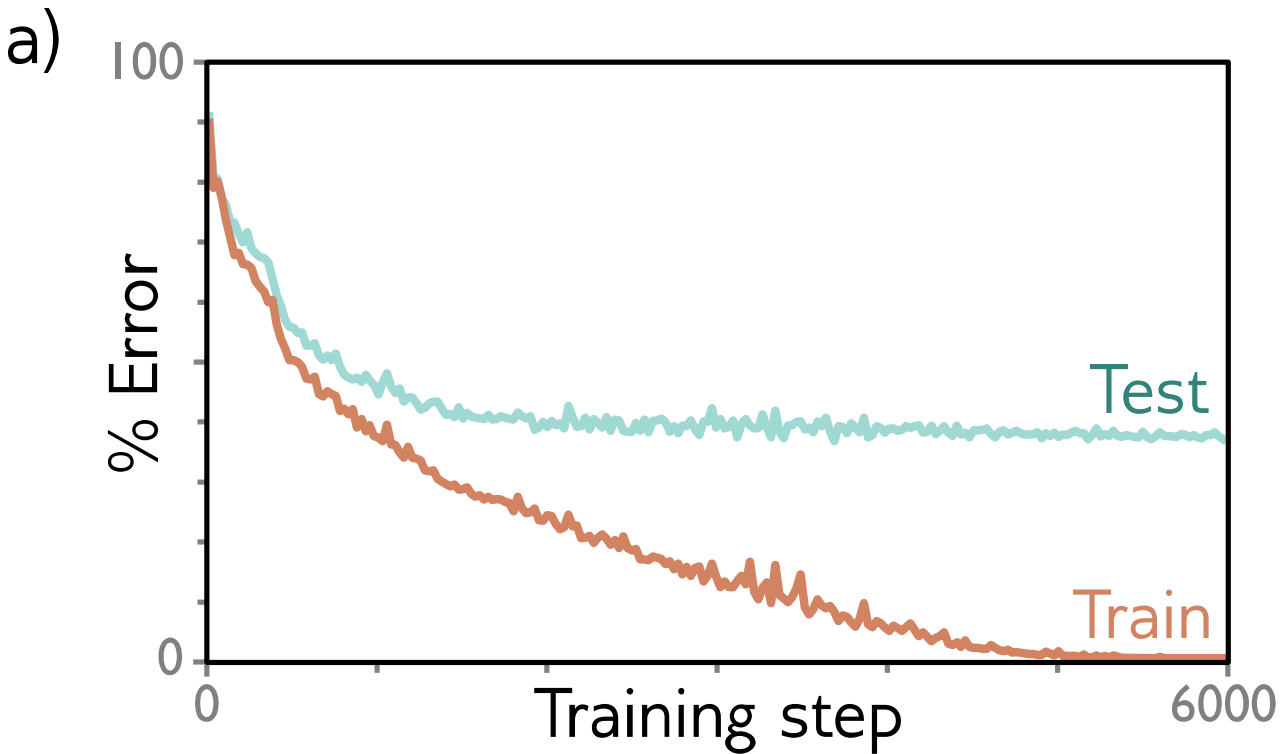
Input size (MB): 0.00  
Forward/backward pass size (MB): 0.00  
Params size (MB): 0.06  
Estimated Total Size (MB): 0.06

```
=====
```

# Results



# Need to use separate test data



The model has not **generalized** well to the new data

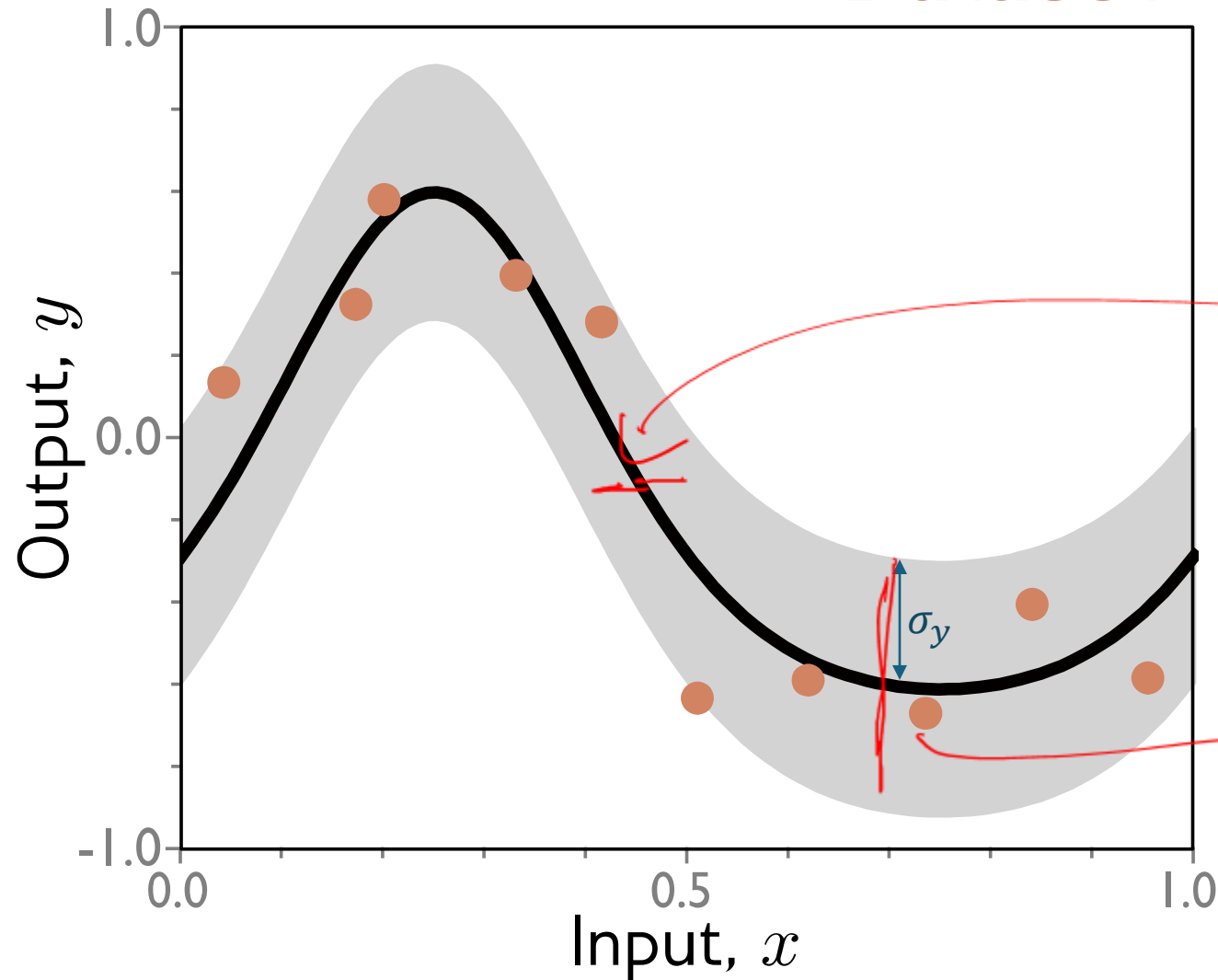
# Any Questions?



- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Choosing hyperparameters

# Regression example with Toy Model

## Dataset



"True" function:

$$y = e^{\sin(2\pi x)}$$

*judge  
✓ this*

Add small uniform noise to  $x$ :

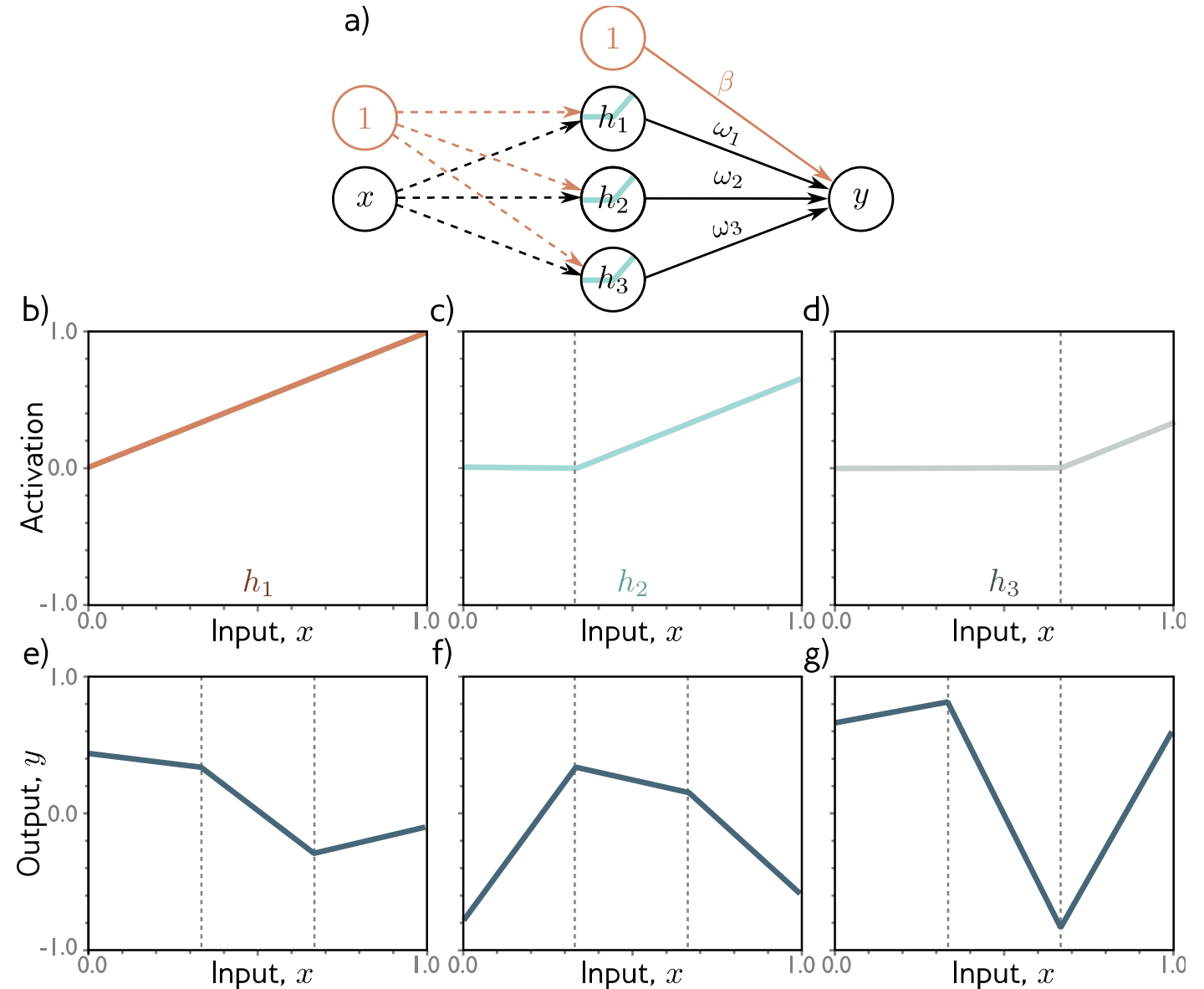
$$x = x + \mathcal{U}(\pm 1/\text{num\_data})$$

Add small Gaussian noise to  $y$ :

$$y = y + \mathcal{N}(0, \sigma_y)$$

# Toy model

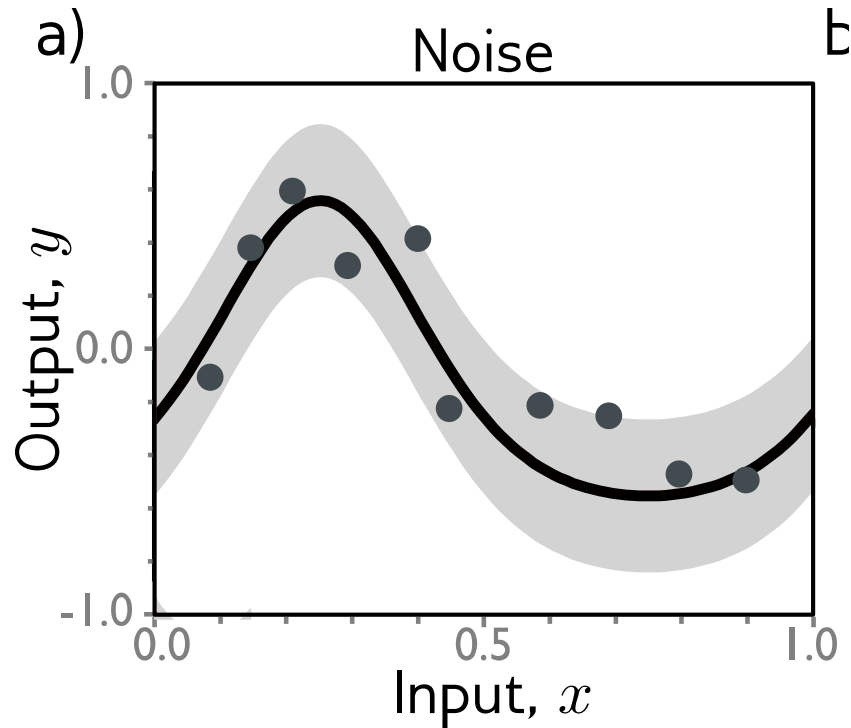
- D hidden units
- First layer fixed so “joints” divide interval evenly, e.g.  $0, 1/D, 2/D, \dots, (D-1)/D$
- Second layer trained
- But... now linear in  $\mathbf{h}$ 
  - so convex cost function
  - can find best solution in closed-form
- A piecewise linear model with D regions.



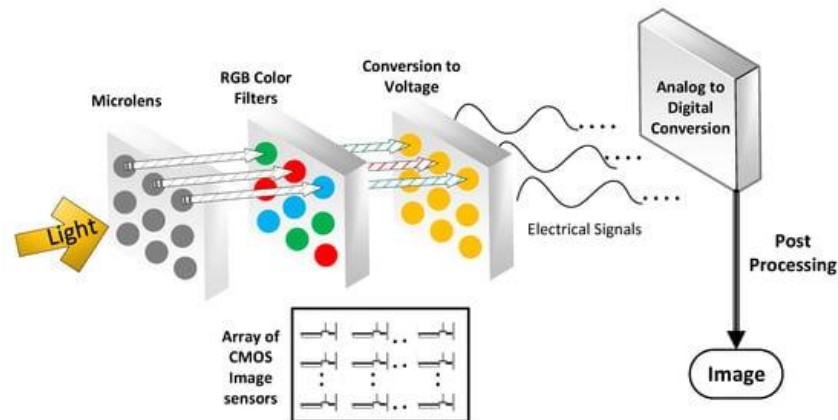
Three possible sources of error:  
*noise*, *bias* and *variance*



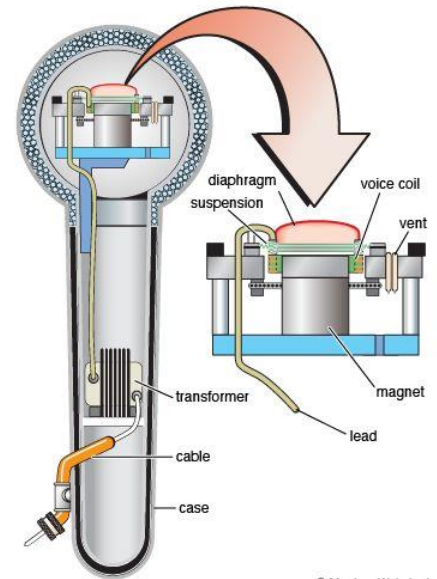
# Noise, bias, and variance



- Genuine stochastic nature of the underlying model
- Noise in measurements, e.g. from sensors
- Some variables not observed
- Data mislabeled



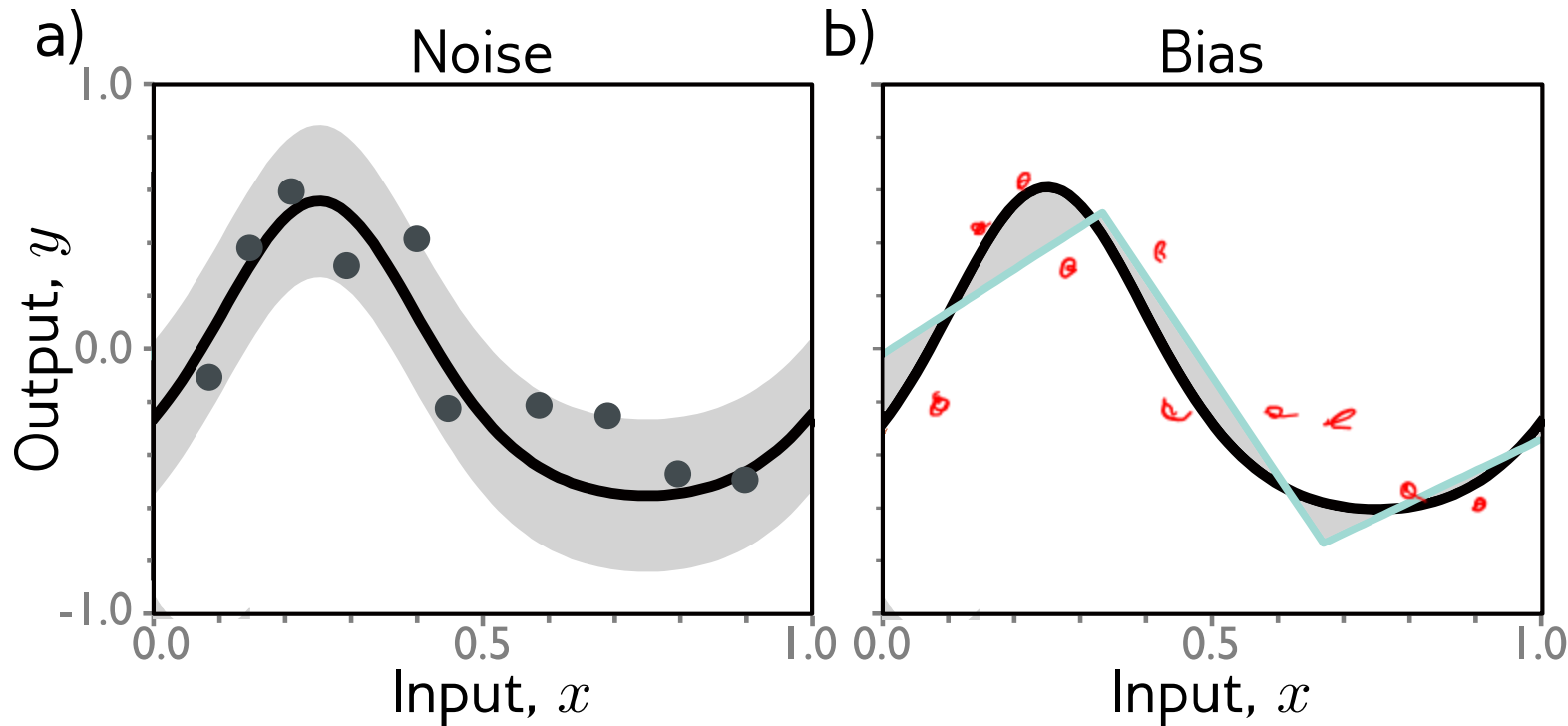
<https://images.app.goo.gl/2PuBhaFpfdL9Pyjb8>



© Merriam-Webster Inc.

<https://images.app.goo.gl/CMDaXSCdX4pqN8Yx7>

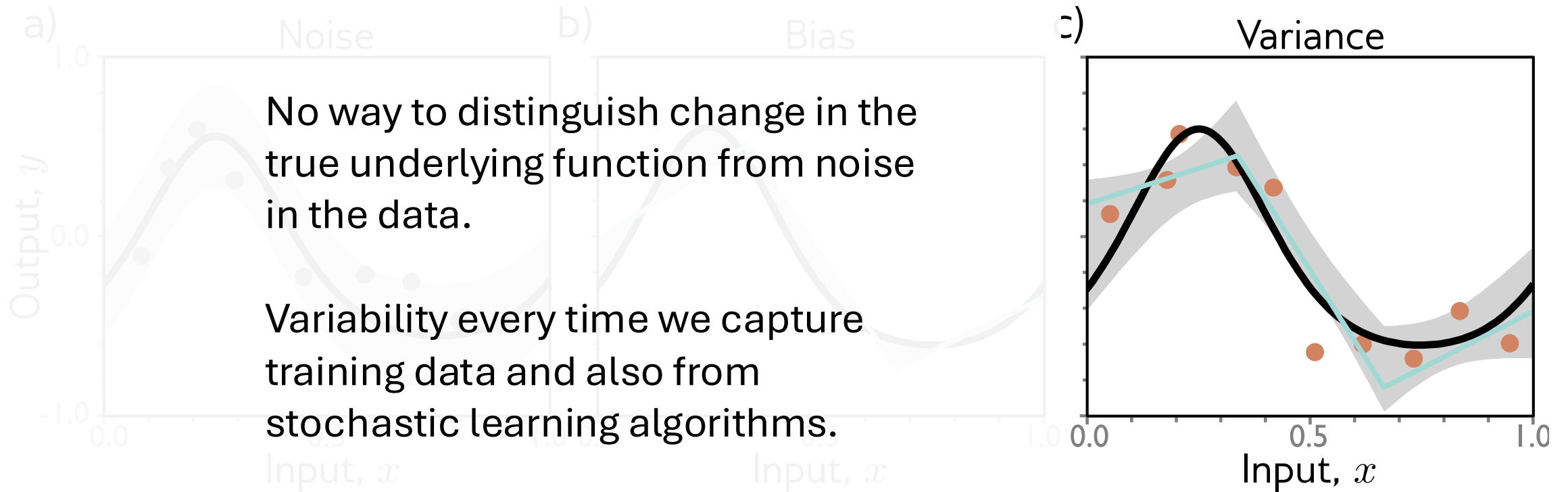
# Noise, **bias**, and variance



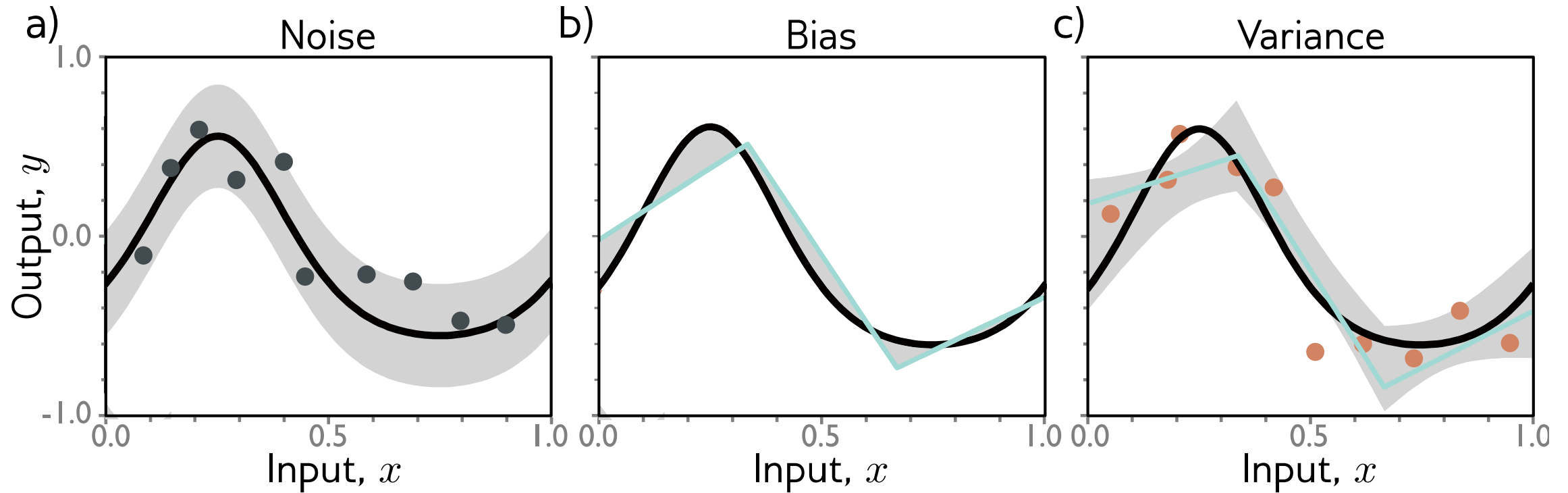
**Bias** occurs because the model lacks precision or capacity to accurately match the underlying function.

E.g. optimal fit with 3 hidden units and 3 line segments

# Noise, bias, and variance



# Noise, bias, and variance



# Least squares regression only

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

- We can show that:

$$\mathbb{E}_y[L[x]] = (f[x, \phi] - \mu[x])^2 + \sigma^2$$

- And then:

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

Expectation over noise in training data  
 Expectation over noise in test data  
 Actual model  
 Best possible model if we had infinite data  
 True function

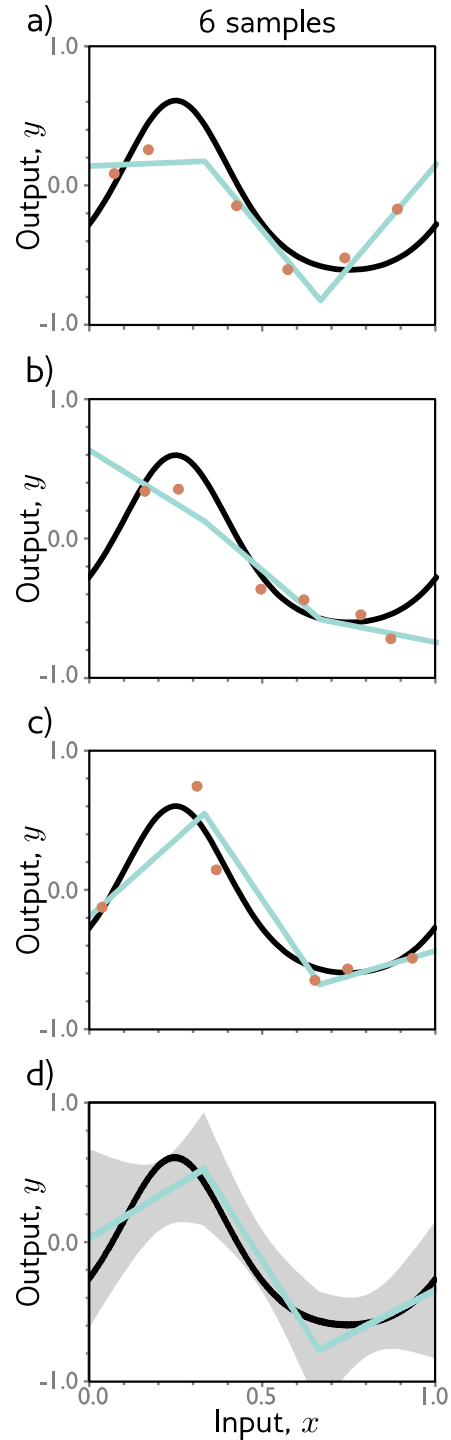
More complex interactions between noise, bias and variance in more complex models.

# Any Questions?



- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Choosing hyperparameters

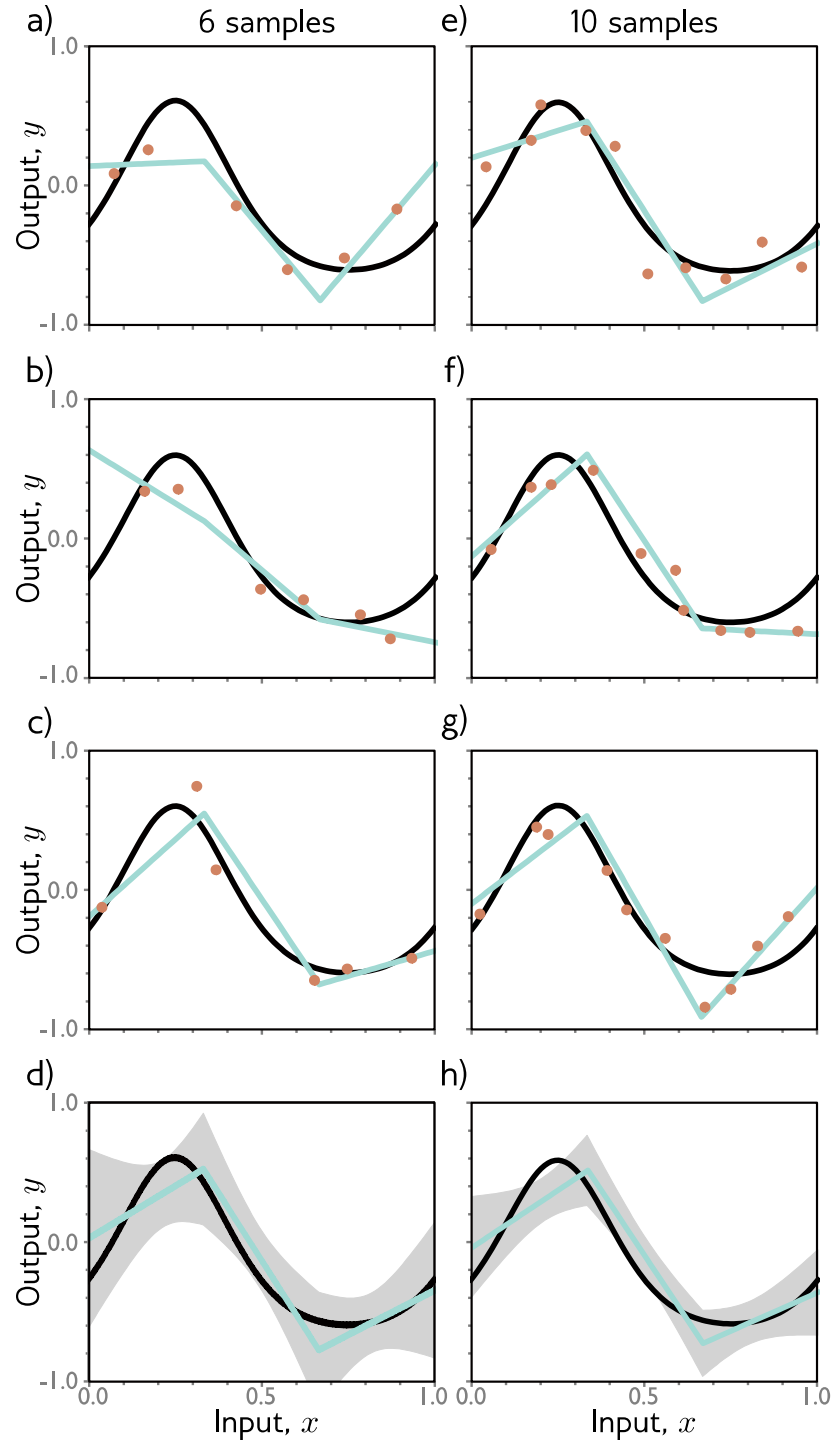
# Variance



When measuring (capturing) 6 different data samples with a fixed model (e.g. 3 hidden units), we get different optimal fits every time.

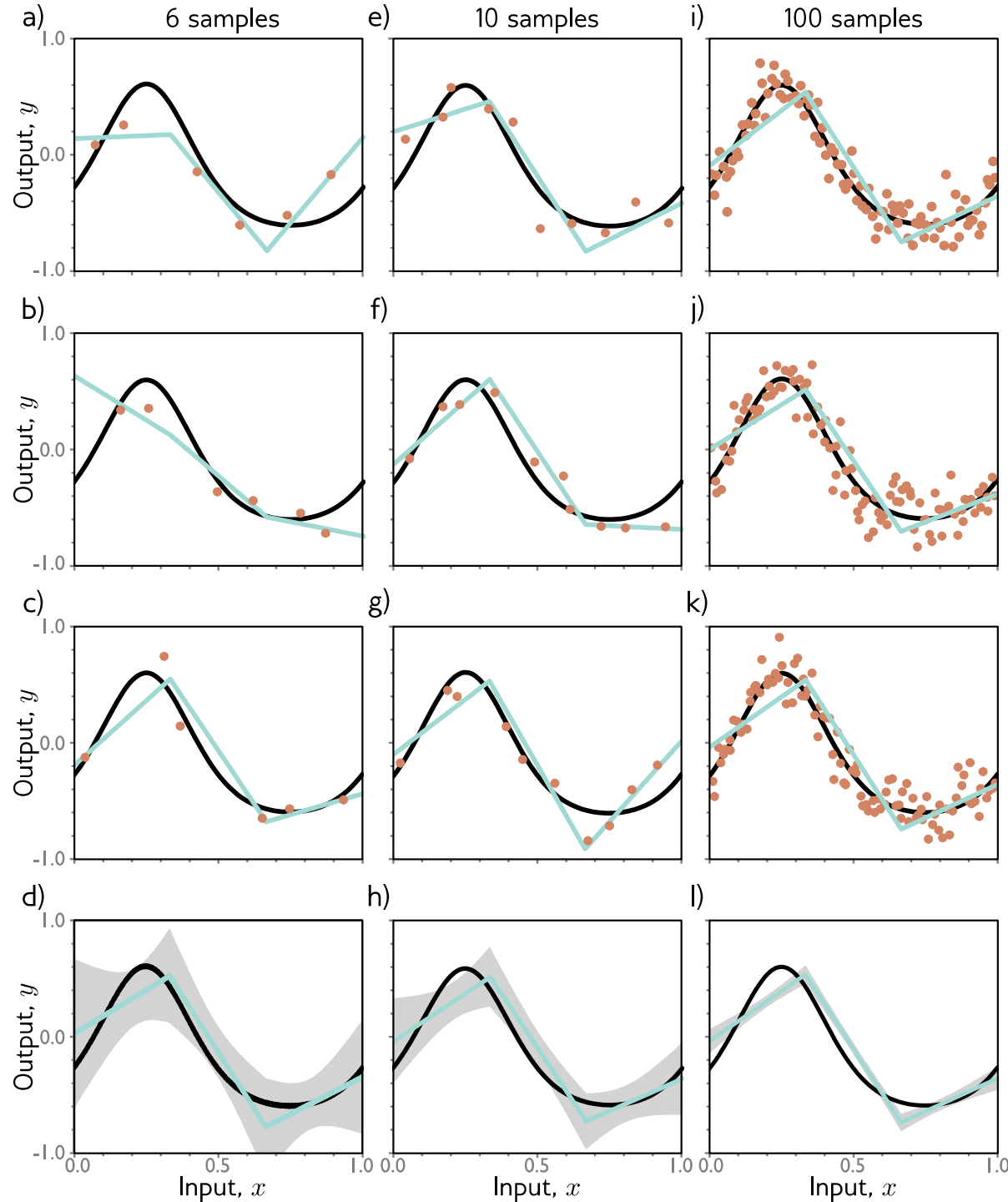


# Variance



Can reduce  
variance by  
adding more  
samples

# Variance



Can reduce  
variance by  
adding more  
samples

Resulting model  
approaches  $f_m$ .  
variance  $\rightarrow \emptyset$ ,  
still have error from  
bias + noise

# Any Questions?

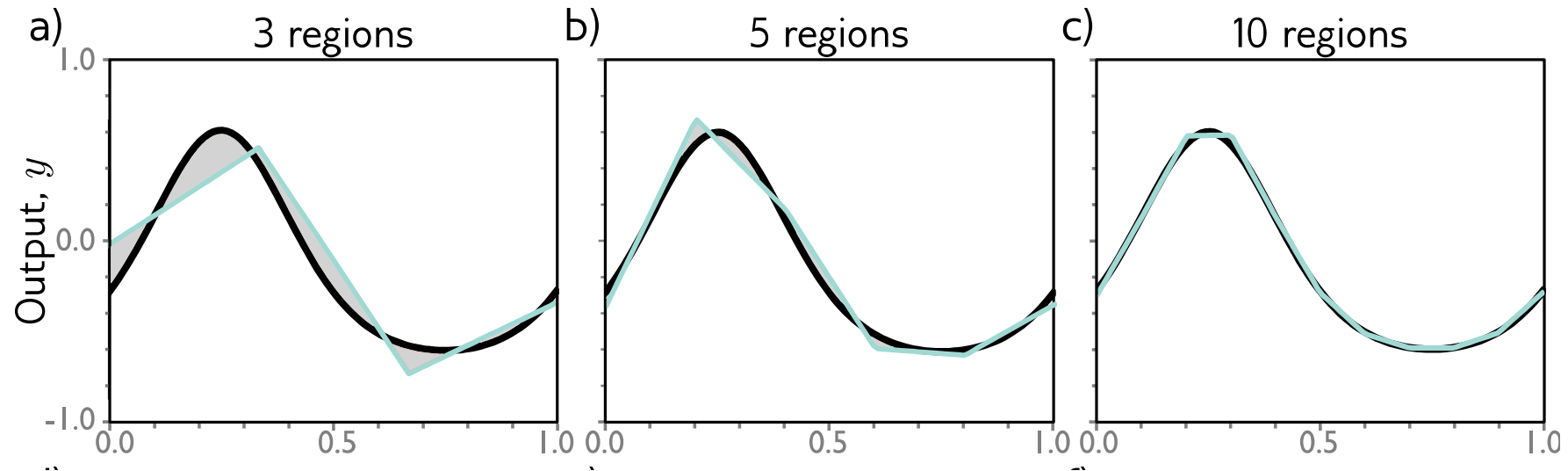


- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Choosing hyperparameters
- Double descent

# Reducing bias

*(example with the true function)*

Bias

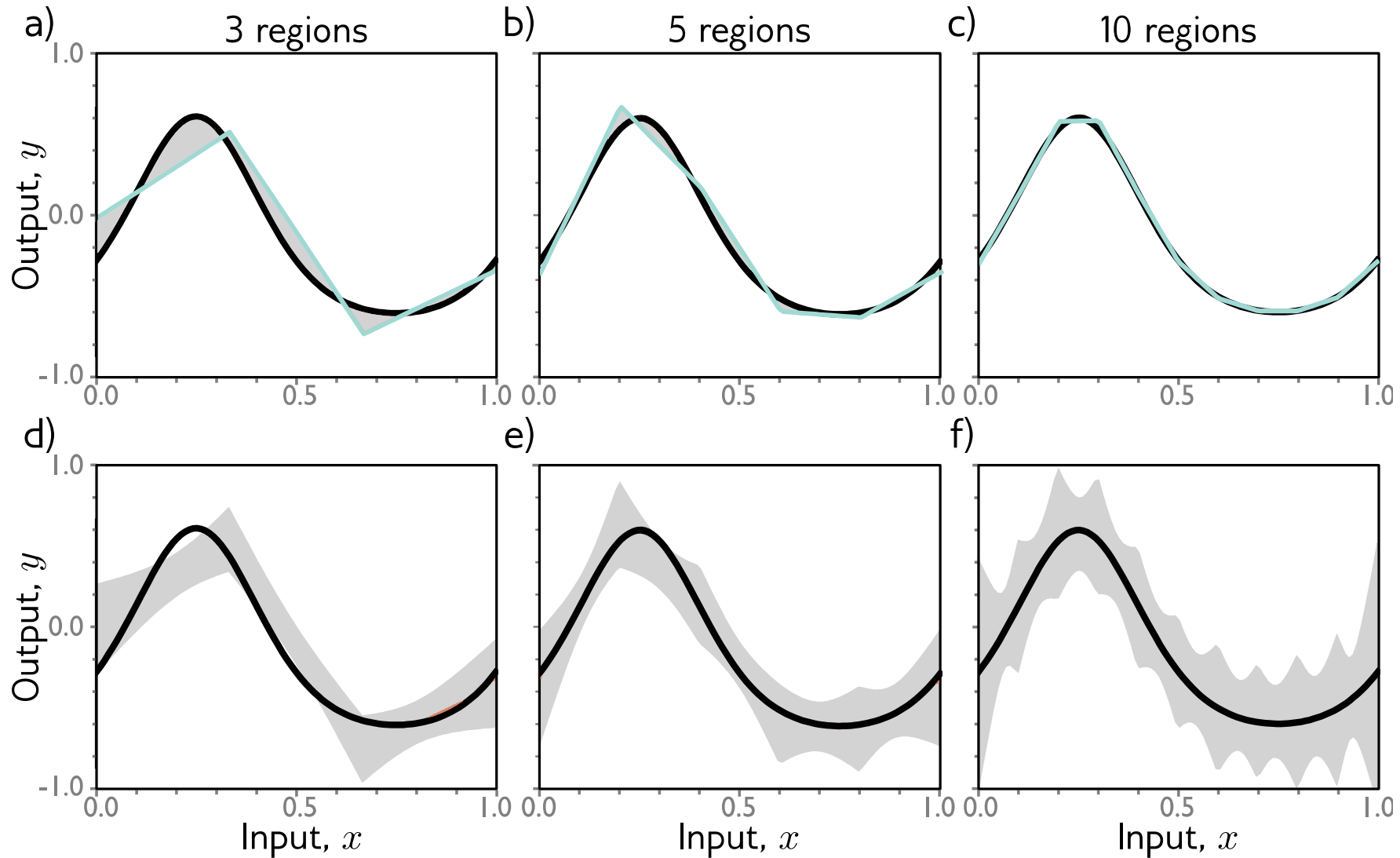


We can reduce bias by adding more model capacity.

In this case, adding more hidden units.

# Reducing bias $\rightarrow$ Increases variance!!

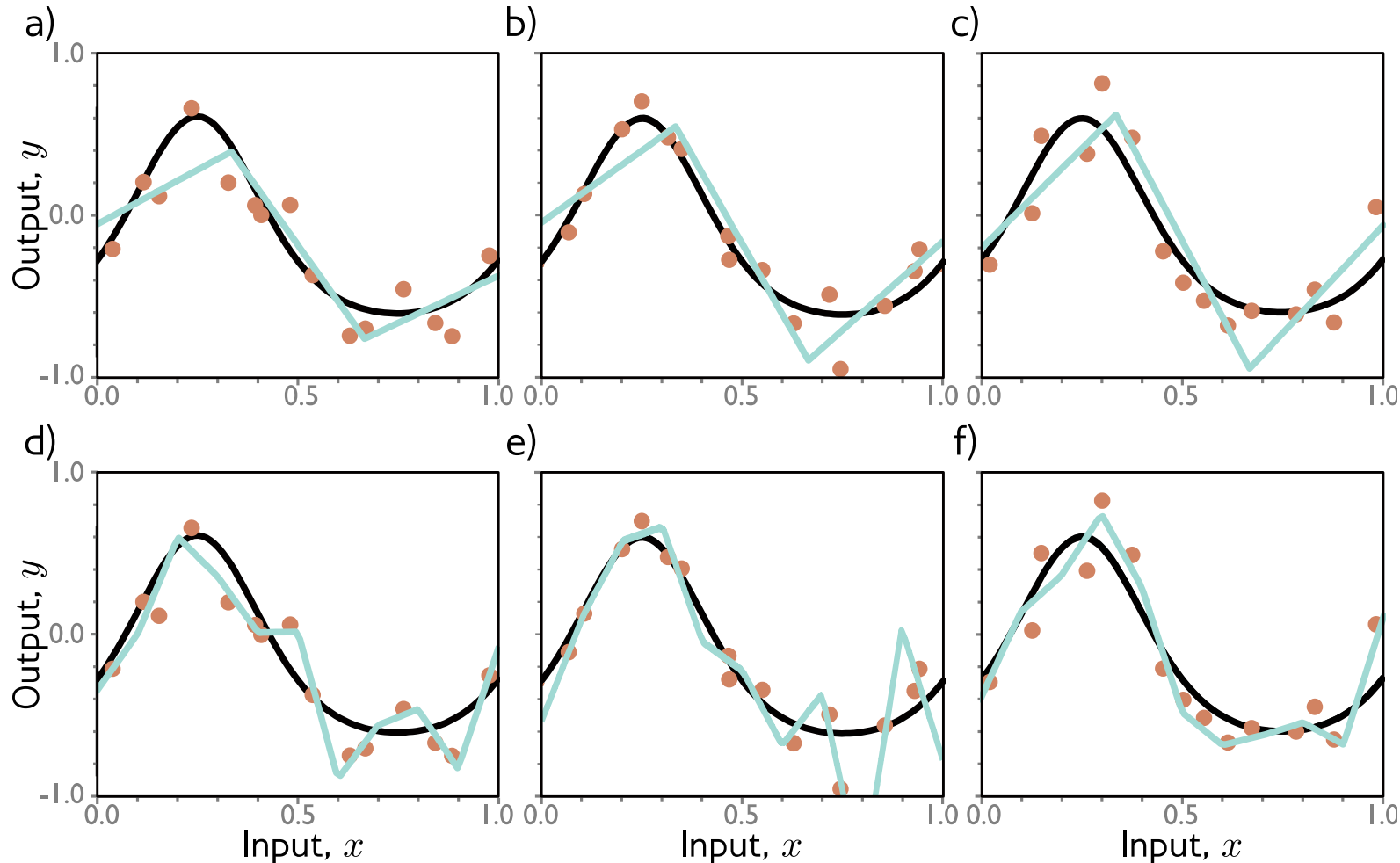
Bias



Variance

# Why does variance increase? Overfitting

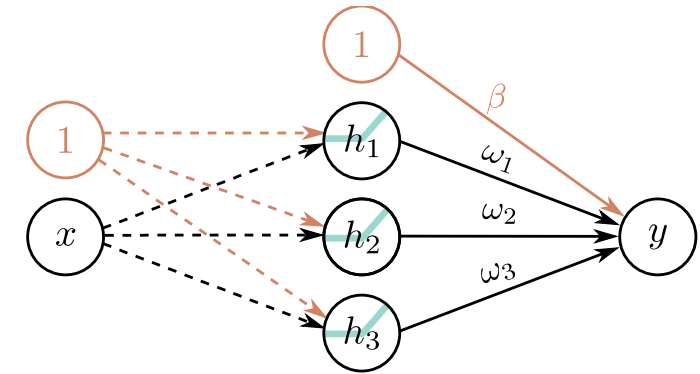
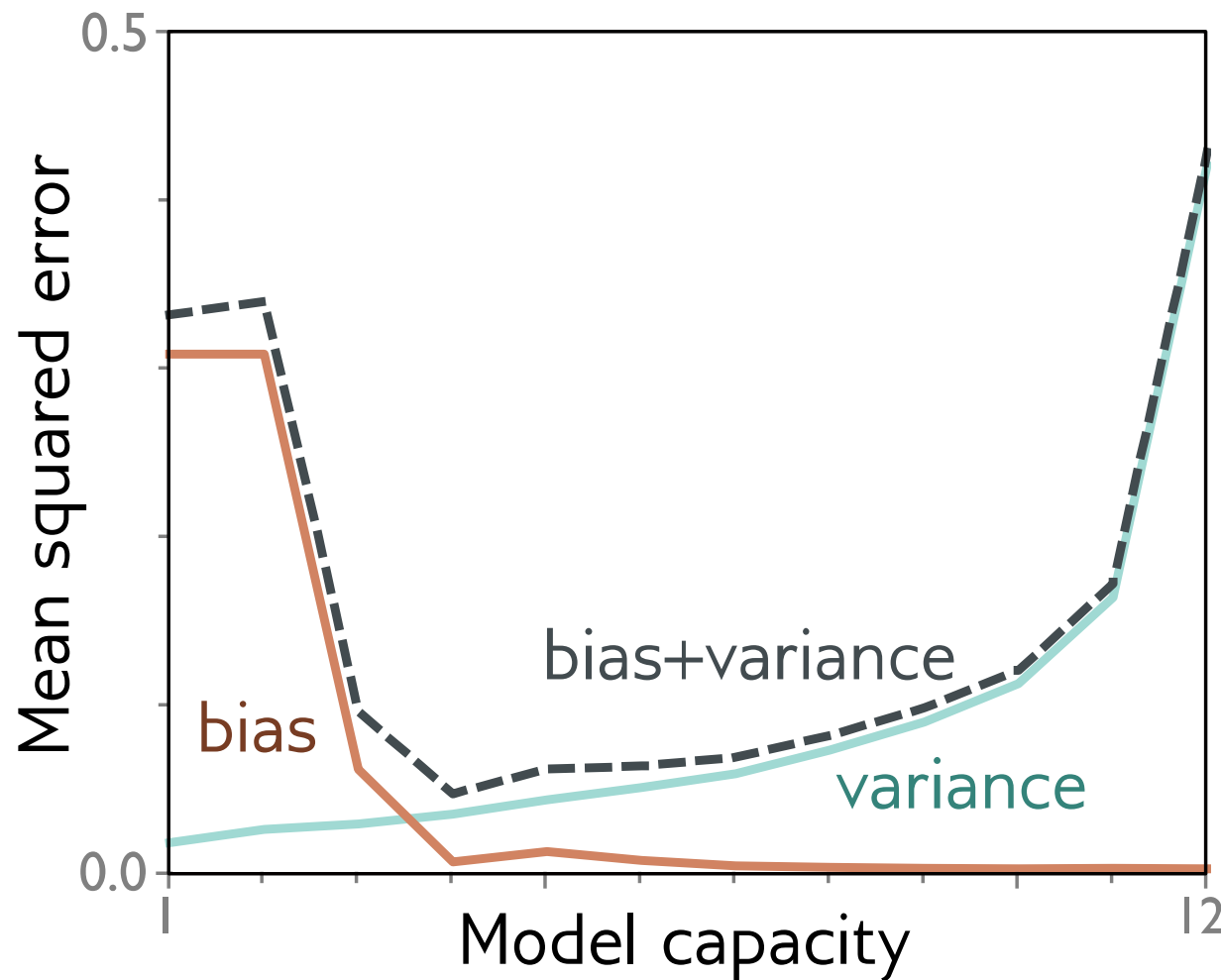
3 Regions



10 Regions

Describes the training data better, but not the true underlying function (black curve)  
Many ways to fit a sample of 15 data points

# Bias and variance trade-off for the simple linear model



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

Number of hidden units



But does picking model capacity to minimize bias & variance hold for more complex data and models?

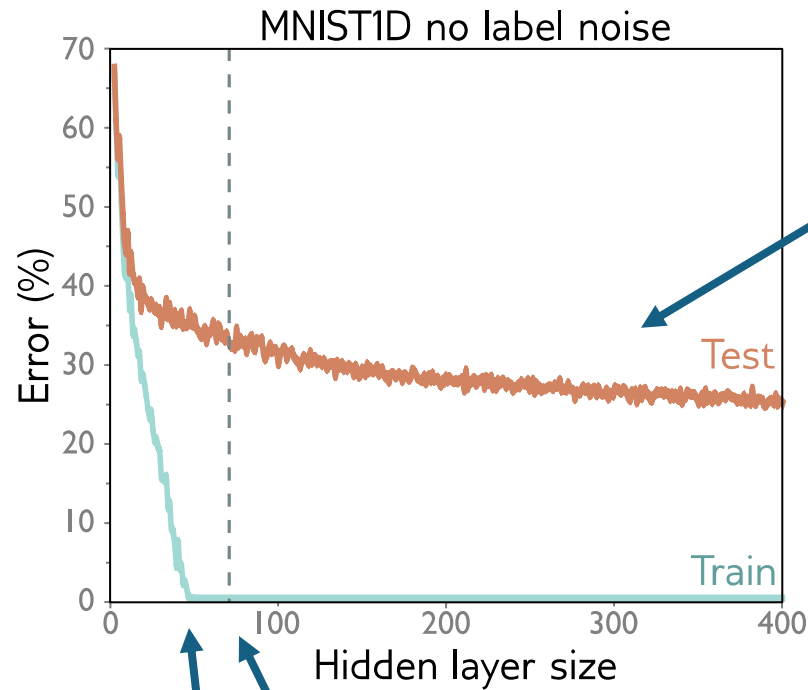
# Any Questions?



- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Choosing hyperparameters

# Train and Test Error versus # of Hidden Layers

- 10,000 training examples
- 5,000 test examples
- Two hidden layers
- Adam optimizer
- Step size of 0.005
- Full batch
- 4000 training steps

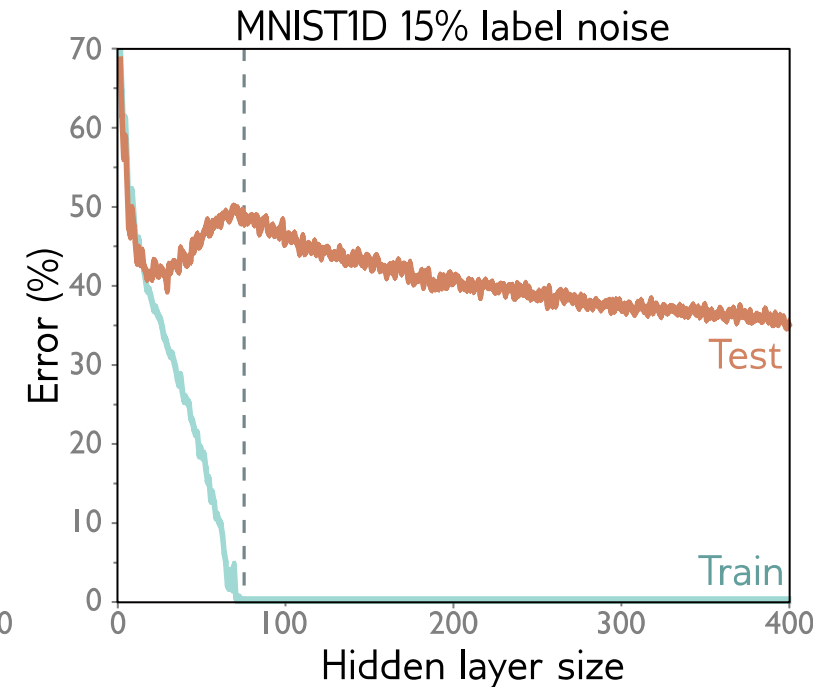
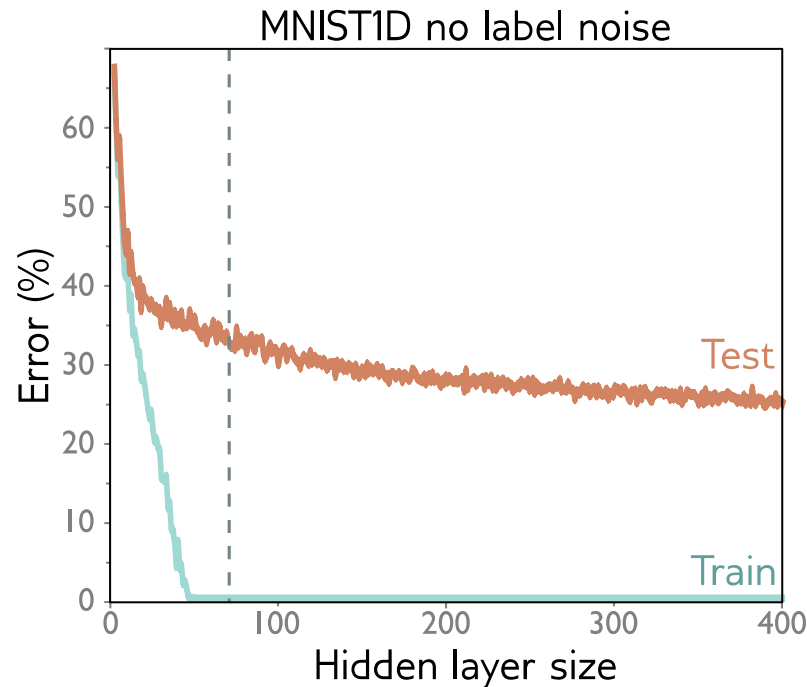


Test error keep decreasing even as we keep increasing model capacity!

Training parameters = Training examples

Model has *memorized* the training set  
Why do we say that?

Now randomize  
15% of the  
training labels

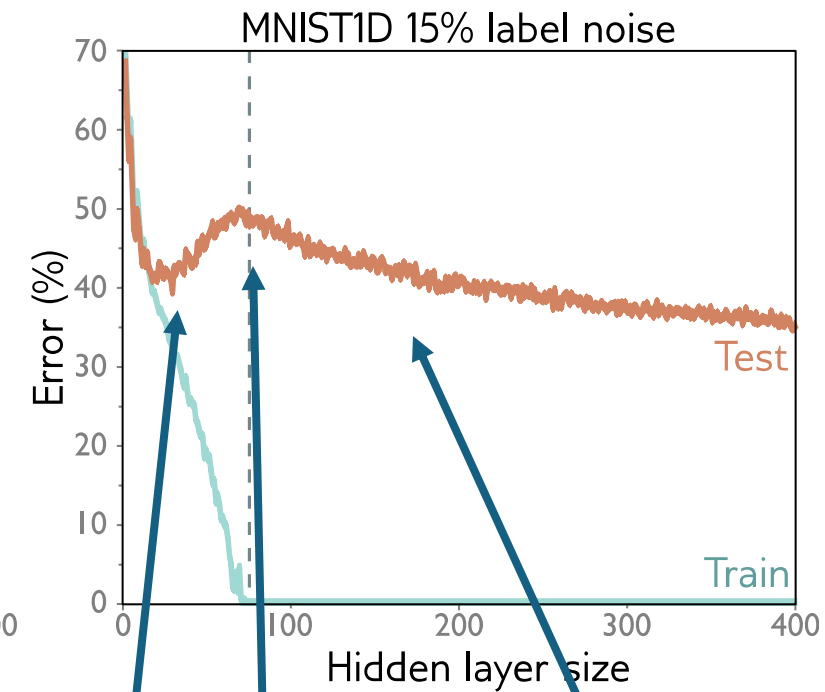
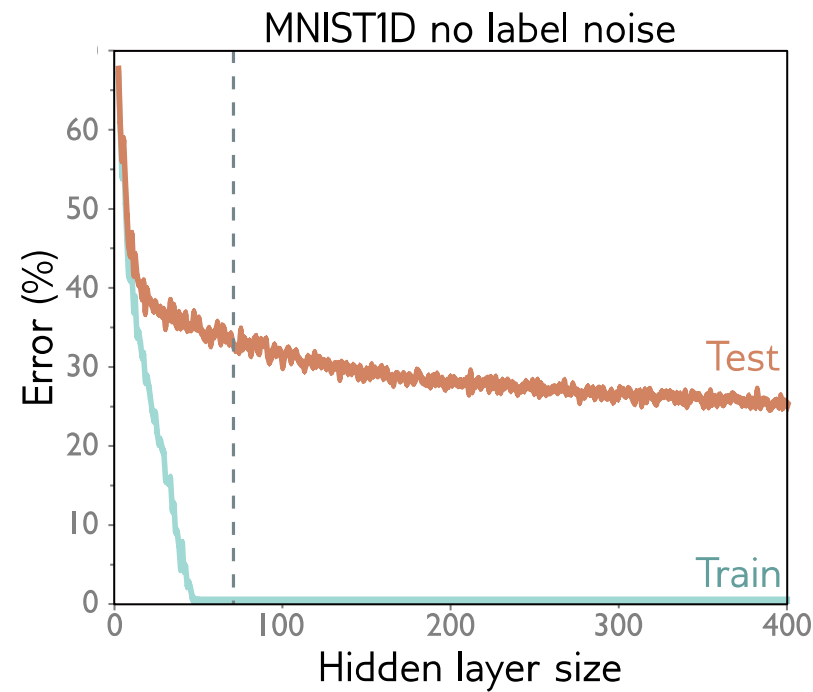


Now we see what looks like bias-  
variance trade-off as we increase  
capacity to the point where the model  
fits training data.

Reminder: vertical dashed line is where:  
# training parameters = # training samples

But then???

# Double Descent



Classical or under-  
parameterized  
regime

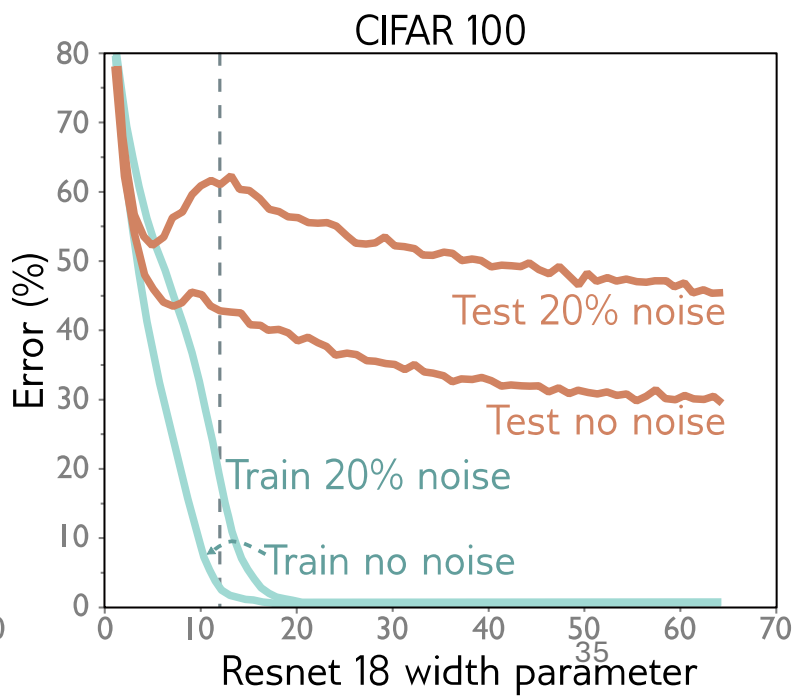
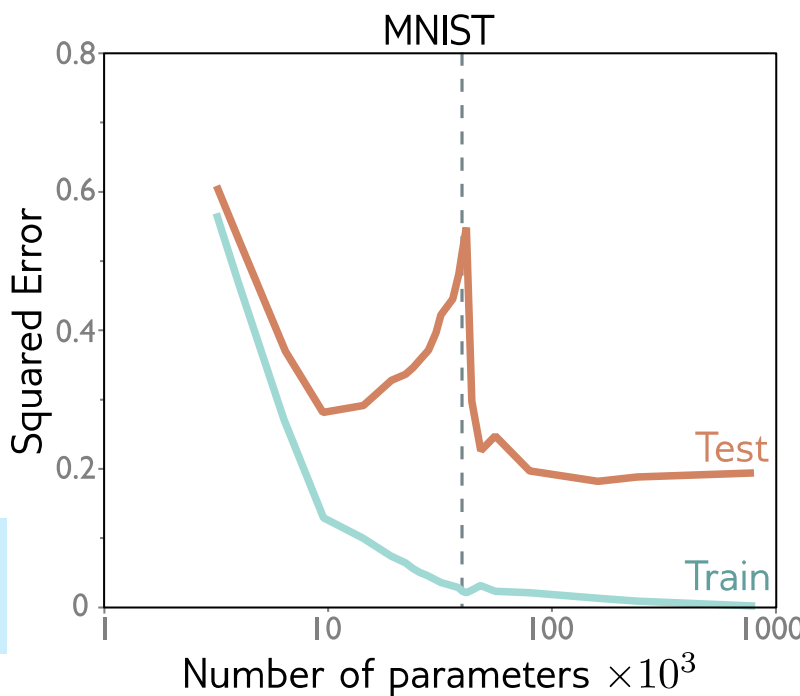
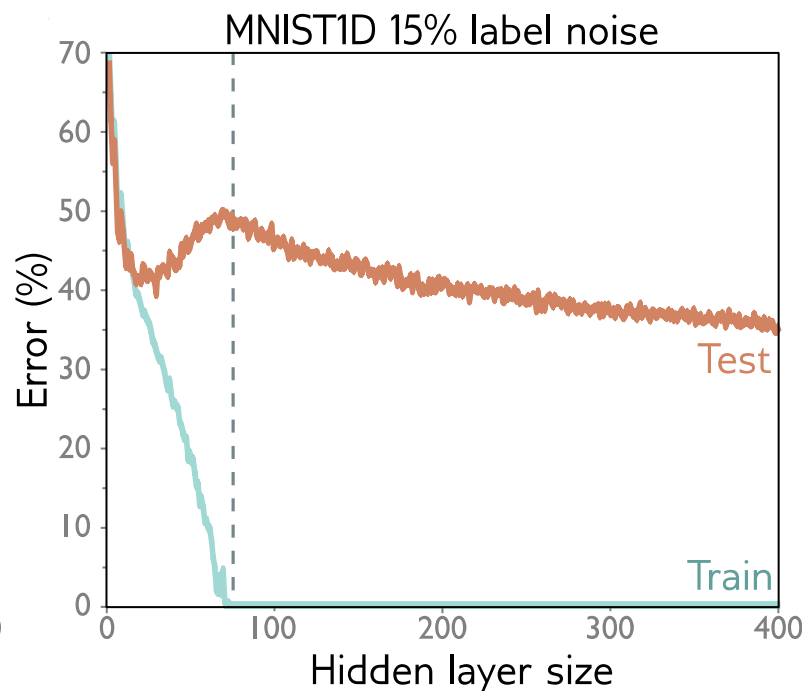
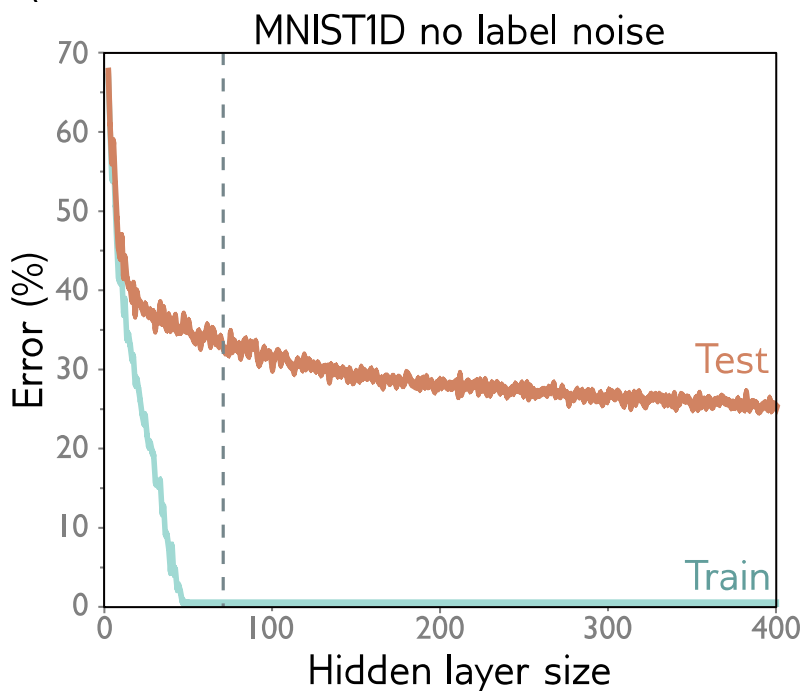
Modern or over-  
parameterized  
regime

Critical regime

Reminder: vertical dashed line is where:  
# training parameters = # training samples

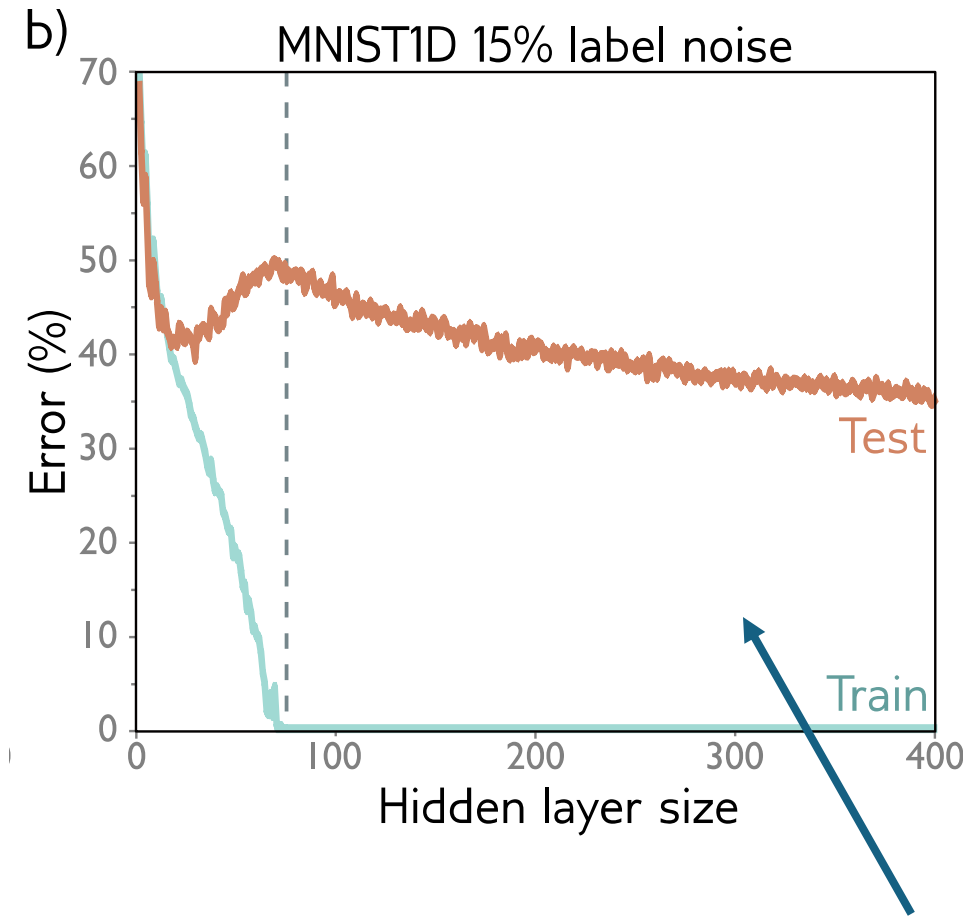
Same  
phenomenon  
shows up on  
MNIST and  
CIFAR100

Reminder: vertical dashed line is where:  
# training parameters = # training samples

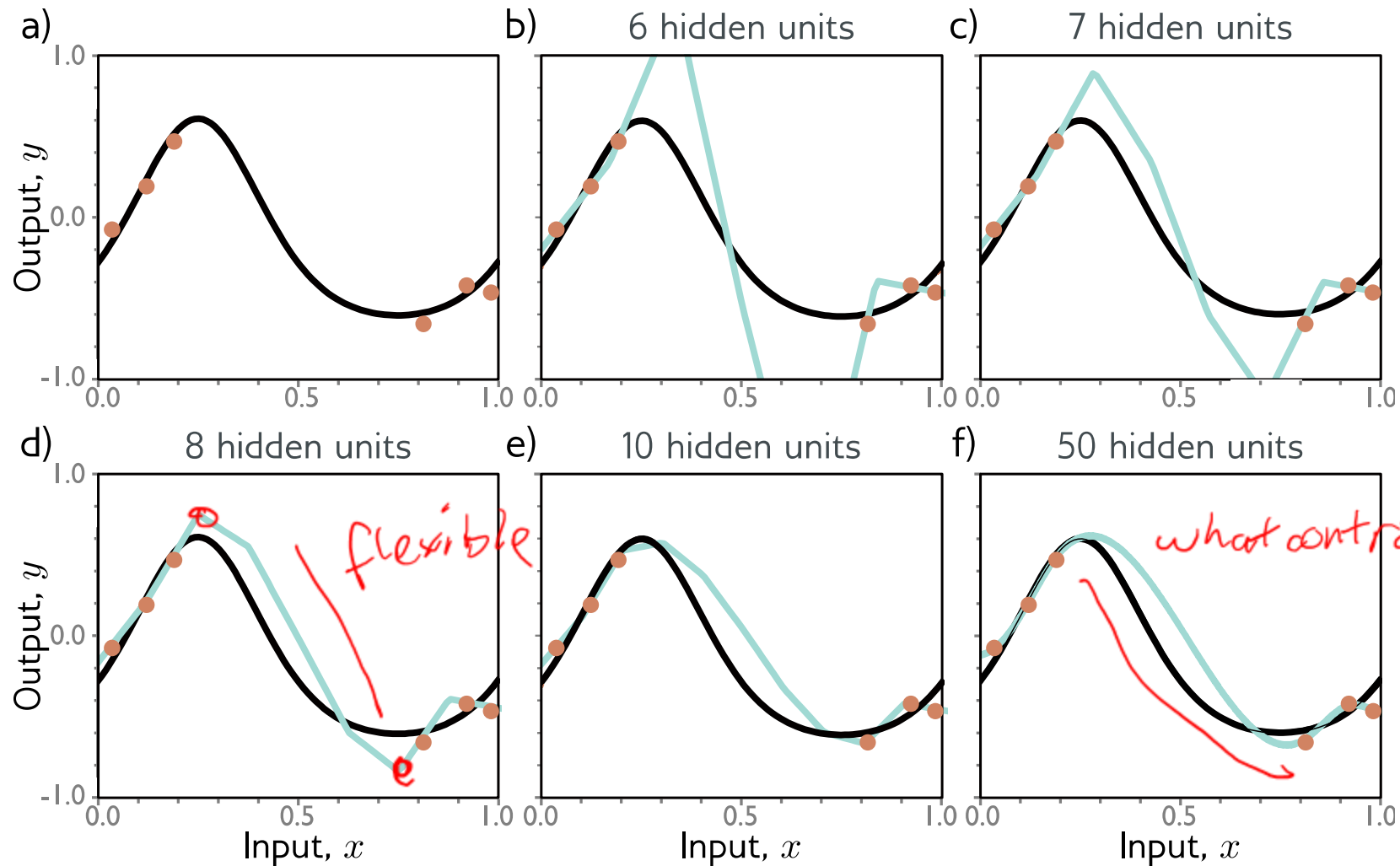


# Double Descent

- Note that training loss is very close to zero.
- Whatever is happening isn't happening at training data points
- Model never sees test set during training
- Must be happening between the data points??





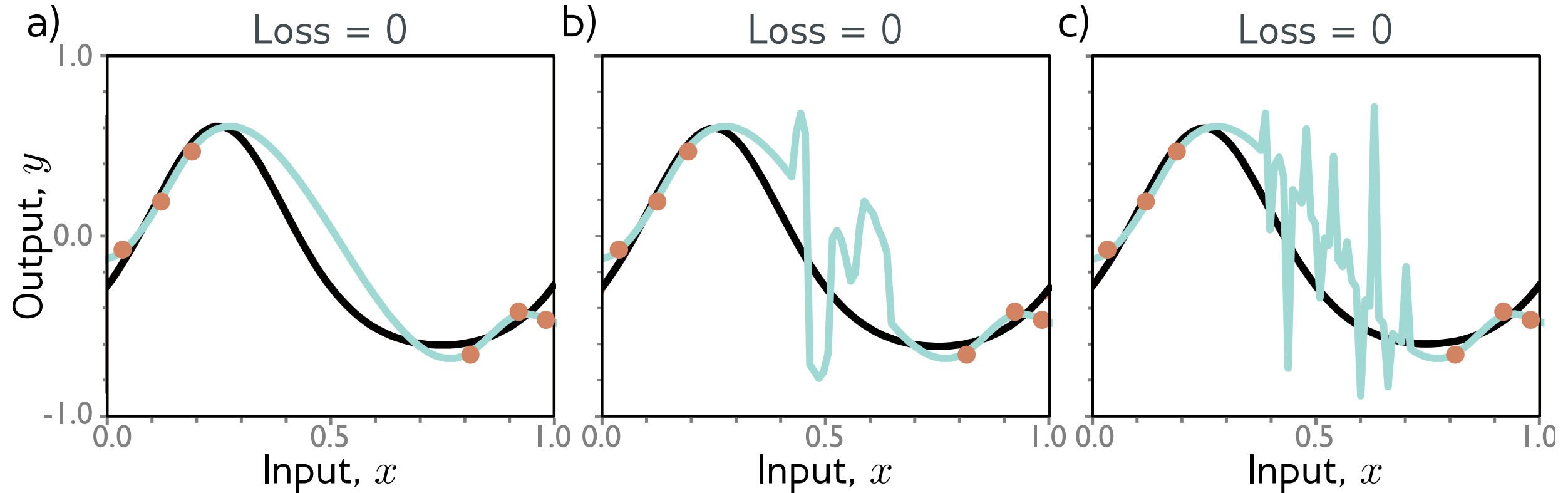


Potential explanation:

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

But why?

# Next Week: How to bias for smoothness?



- 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**

# Any Questions?



- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- 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

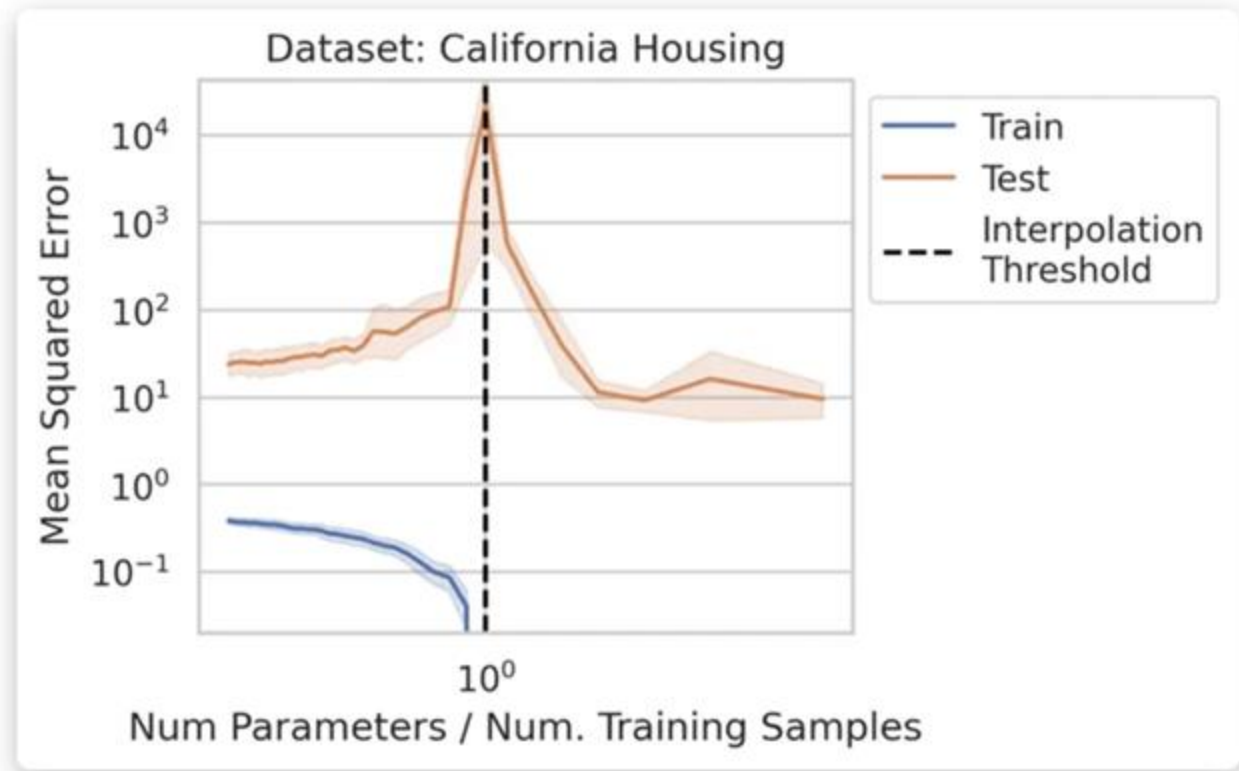
# Any Questions?



- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Choosing hyperparameters

# Double Descent Demystified

- Linear regression too?

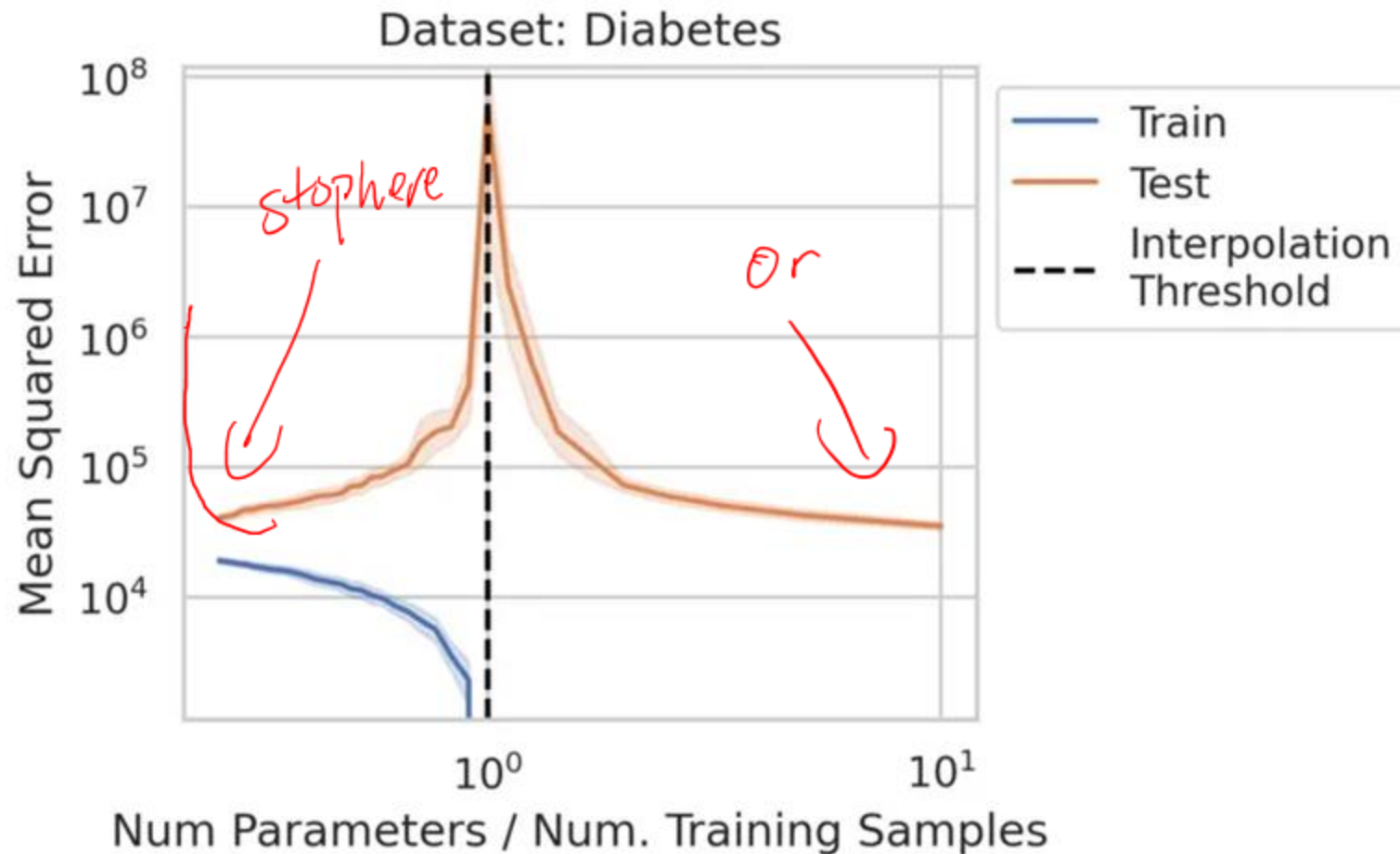


<https://iclr-blogposts.github.io/2024/blog/double-descent-demystified/>

Manufactured polynomial columns from original columns

$x^2y^3$

# Reproduced with Classic Data Sets



# Double Descent Test Setup

- Used polynomial features to generate arbitrary numbers of features
- When more polynomial features than training samples, regression has multiple parameters for exact fit.
  - Pick parameters minimizing a norm ← this is a regularization.



# Rough Explanation for Double Descent

## 1. Parameters $\ll$ samples

- Model can only fit overall trends. Cannot fit individual points particularly well.
- Training and test loss improve with more parameters.

## 2. Parameters $\sim$ samples $\leftarrow$ overfitting b/c trying to interpolate.

- Model can barely / not quite fit all training points.
- Contortions are likely.
- Detailed analysis says “singular values” a lot. Read the paper if curious.

## 3. Parameters $\gg$ samples

- Model can easily fit all training points.
- Lots of freedom to make parameter norms smaller.
- Some intuitive and proven connections to better generalization from smaller norms.

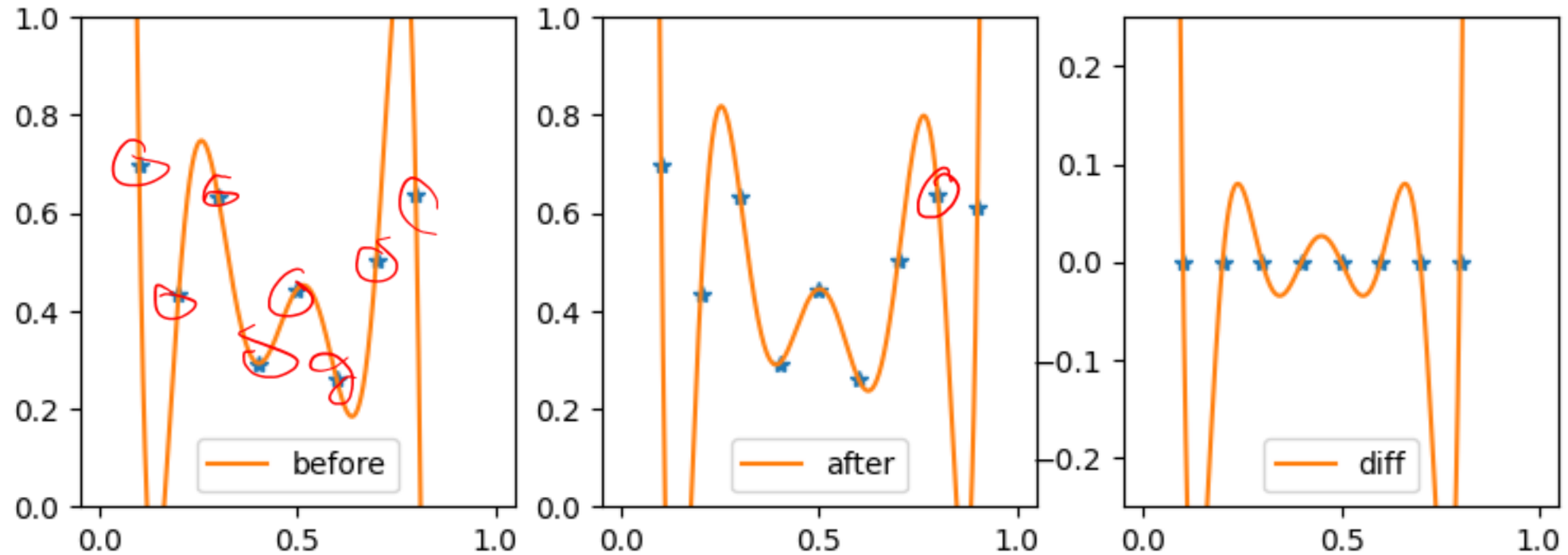
$$y = e^x$$

$$\text{If } y = 3x_1^2 + 2x_1^3 x_2^4$$

cols  $x_1, x_2, x_1^2, x_1 x_2, x_2^2$

# Re: Contortions are Likely

Original slide title: What Happens When You Add Another Point?



# Any Questions?

???

# Theoretical Results

generalization is opposed to fast output changes w/ small input changes

Takeaways from extra readings? Big weights make later values change faster

- For Valid Generalization the Size of the Weights is More Important than the Size of the Network
- Train faster, generalize better: Stability of stochastic gradient descent

output change  $\propto T$  (size factor per layer)

Watch for these themes as we go through regularization examples.

# Any Questions?

???

# Regularization

- Why is there a generalization gap between training and test data?
  - Overfitting (model describes statistical peculiarities)
  - Model unconstrained in areas where there are no training examples
- **Regularization** = methods to reduce the generalization gap
- Technically means adding terms to loss function
- But colloquially means any method (hack) to reduce gap between training and test data

# Regularization

- Explicit regularization
- Implicit regularization
- Early stopping
- Ensembling
- Dropout
- Adding noise
- Transfer learning, multi-task learning, self-supervised learning
- Data augmentation