

Measuring Performance

DL4DS – Spring 2024

Where we are



=== Foundational Concepts ===

- ✓ 02 -- Supervised learning refresher
 - ✓ 03 -- Shallow networks and their representation capacity
 - ✓ 04 -- Deep networks and depth efficiency
 - ✓ 05 -- Loss function in terms of maximizing likelihoods
 - ✓ 06 -- Fitting models with different optimizers
 - ✓ 07a – Gradients on deep models and backpropagation
 - ✓ 07b – Initialization to avoid vanishing and exploding weights & gradients
-
- 08 – Measuring performance, test sets, overfitting and double descent
 - 09 – Regularization to improve fitting on test sets and unseen data

=== Network Architectures and Applications ===

- 10 – Convolutional Networks
- 11 – Residual Networks
- 12 – Transformers
- Large Language and other Foundational Models
- Generative Models
- Graph Neural Networks
- ...

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

MNIST1D

Scaling down Deep Learning

Sam Greydanus¹

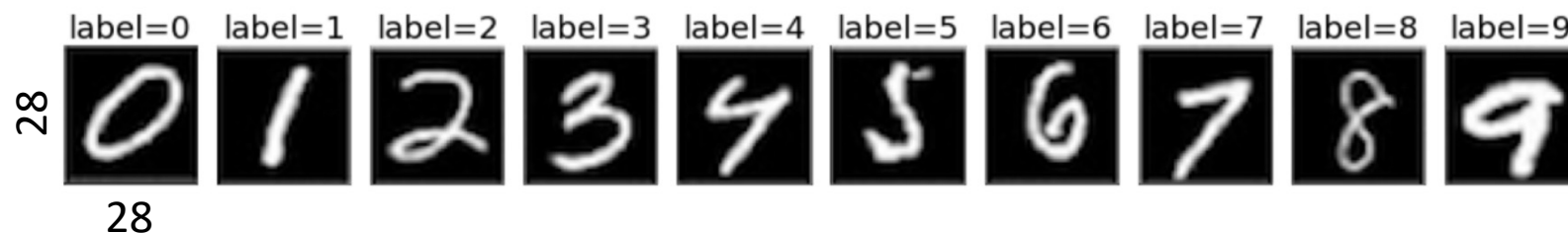
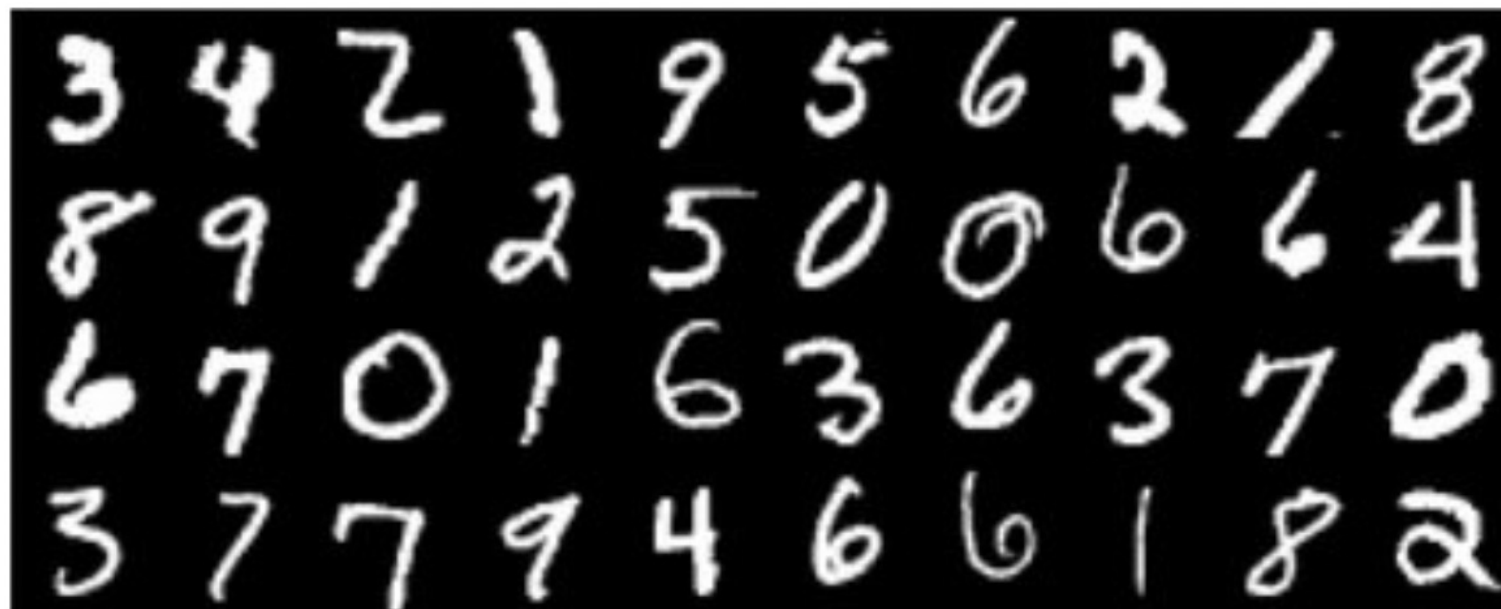
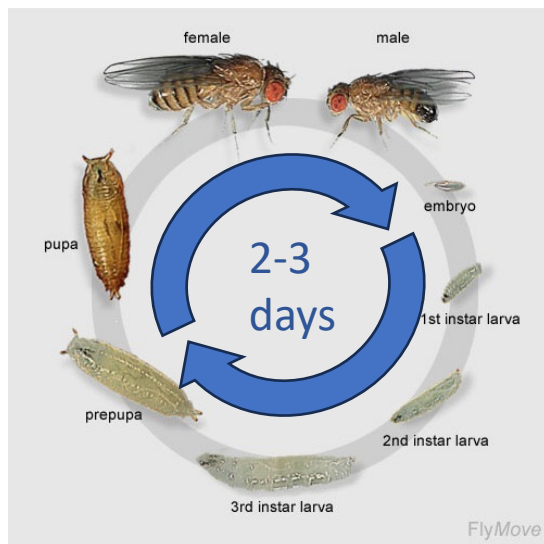
“A large number of deep learning innovations including [dropout](#), [Adam](#), [convolutional networks](#), [generative adversarial networks](#), and [variational autoencoders](#) began life as MNIST experiments. Once these innovations proved themselves on small-scale experiments, scientists found ways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

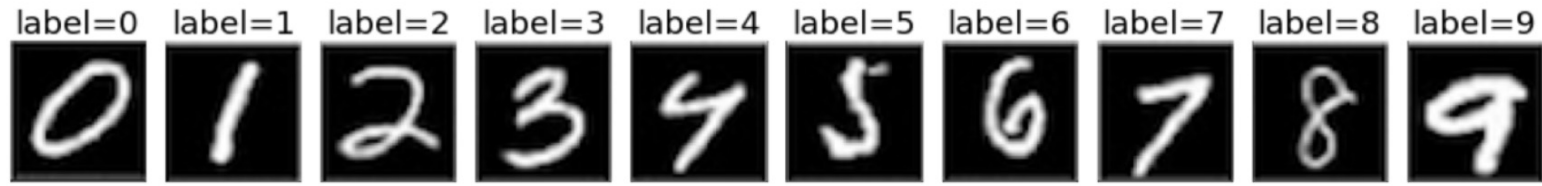


But poorly differentiates model performance:

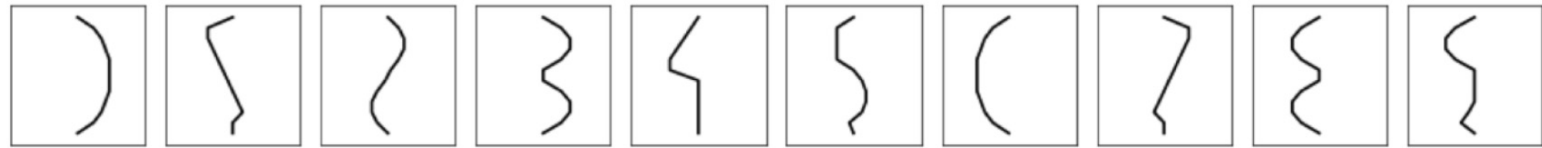
Model Type	Accuracy
Logistic Regression	94%
MLP	99+%
CNN	99+%

MNIST 1D Dataset

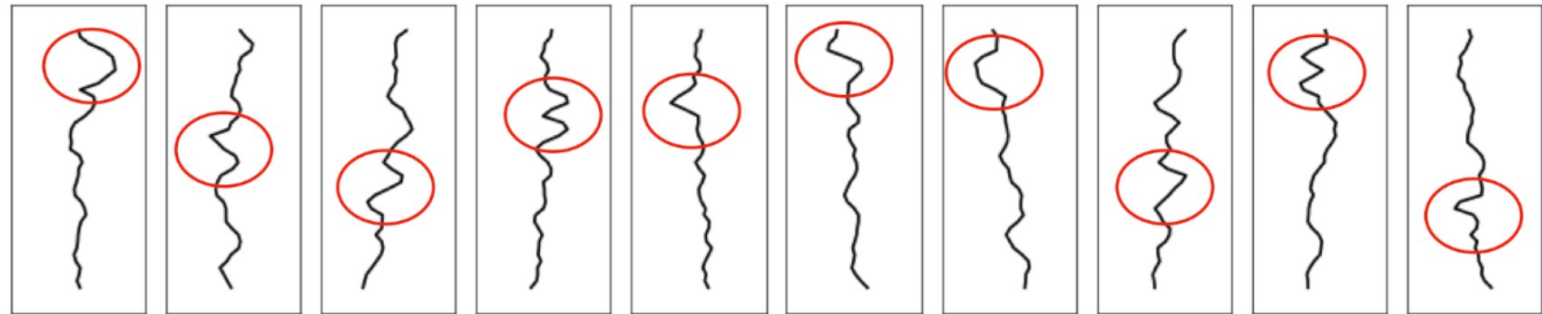
Original MNIST examples



Represent digits as 1D patterns



Pad, translate & transform



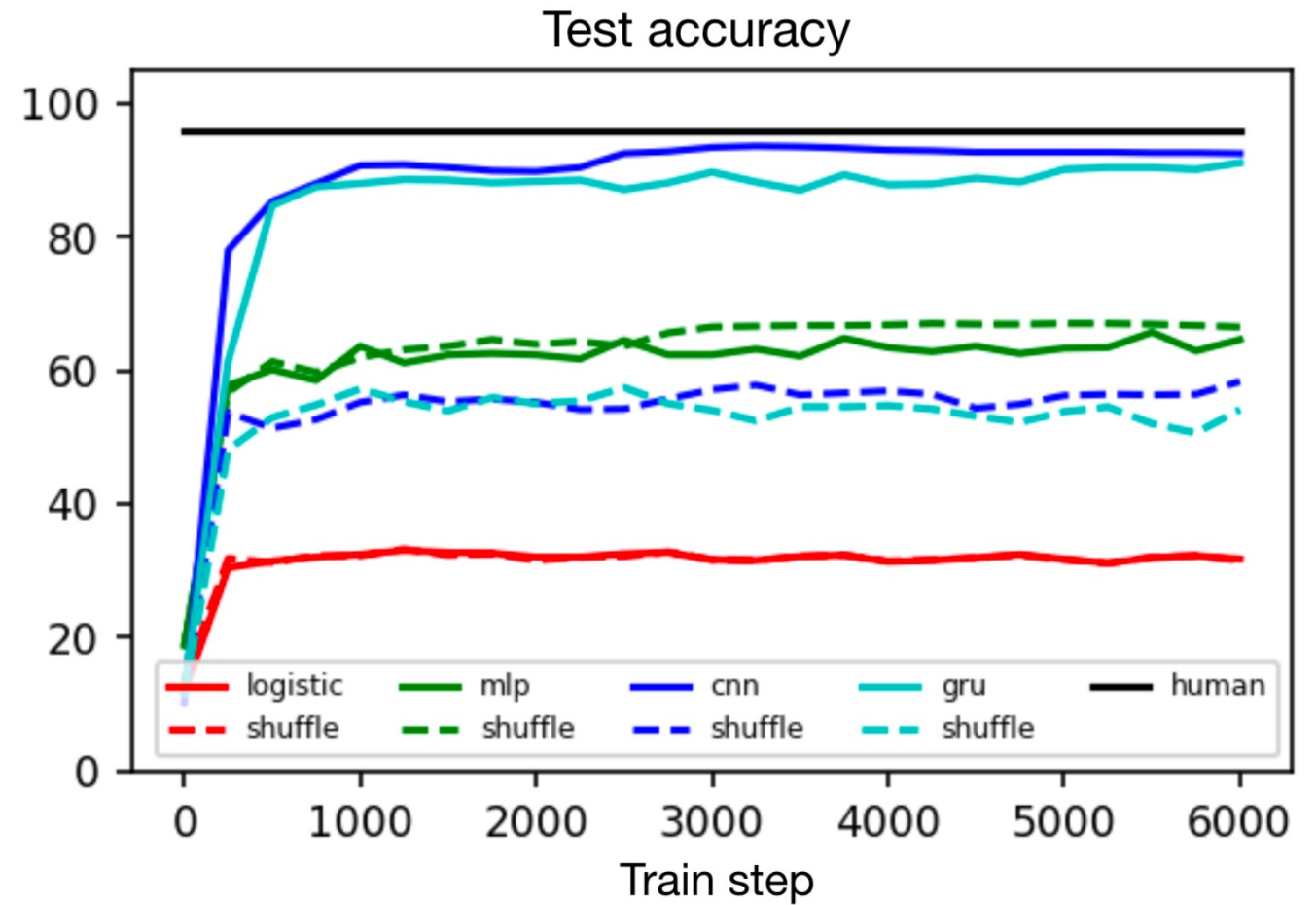
Fixed, 1-D, length-12
templates for each of 10 digit
classes

Generate dataset by
programmatically applying 6
parametric transformations.

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

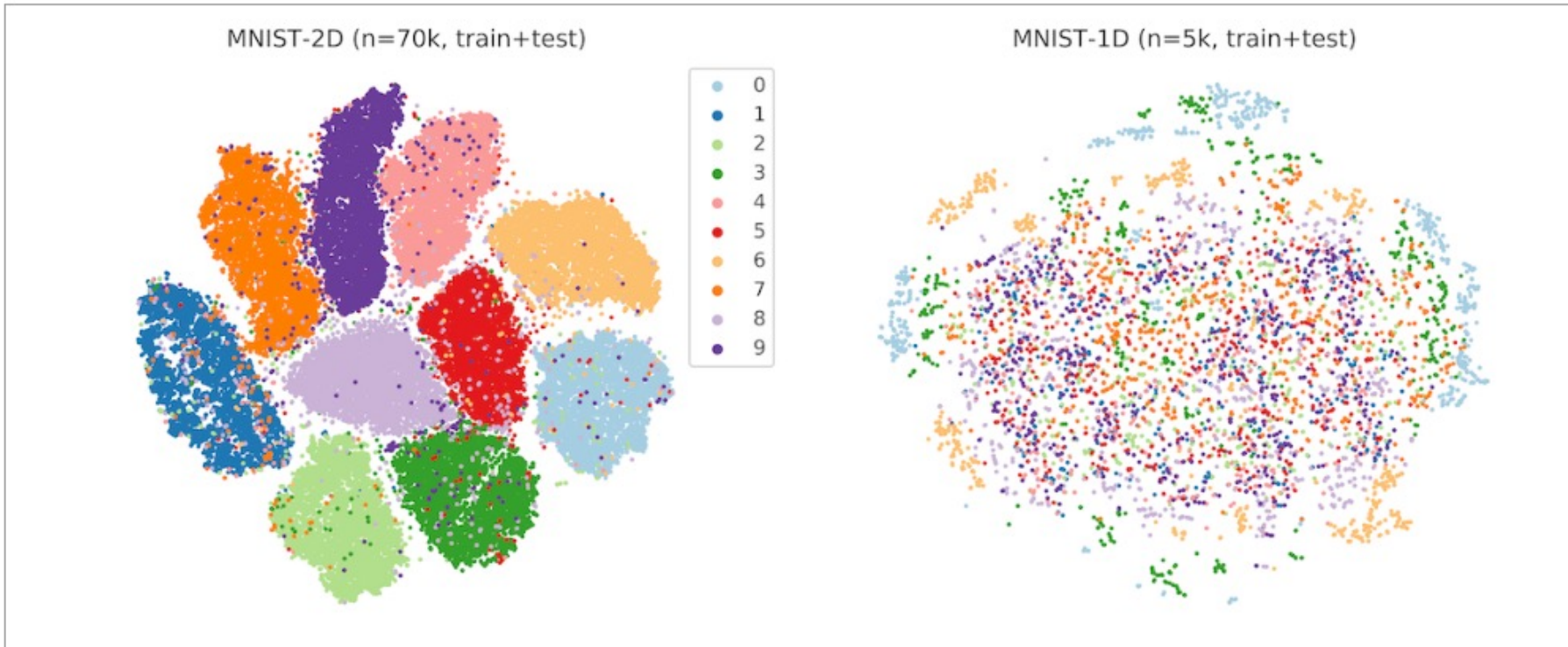
MNIST 1D

Differentiates performance of different model types much more than MNIST



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

Visualizing MNIST and MNIST-1D with tSNE

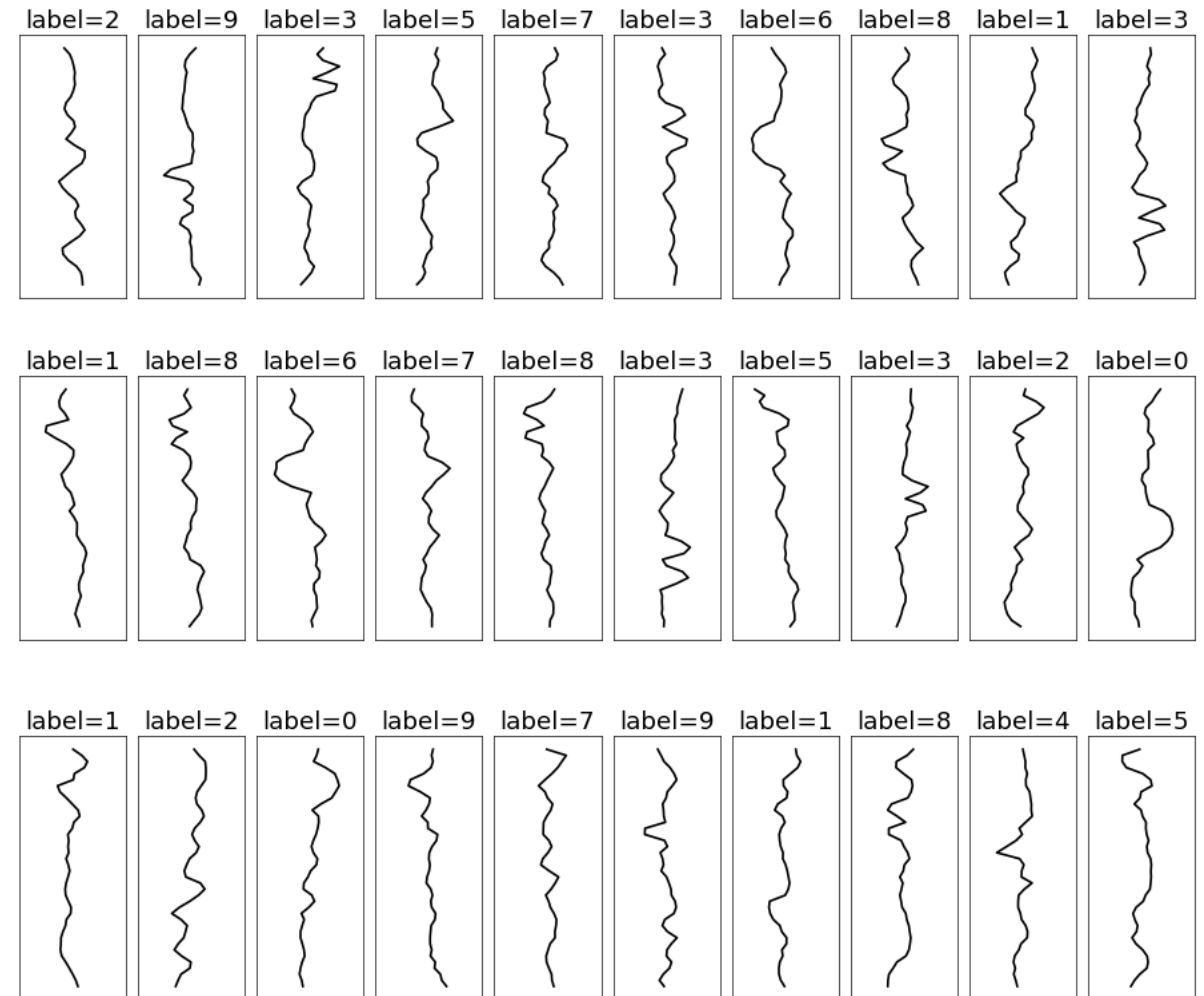


Visualizing the MNIST and MNIST-1D datasets with tSNE. The well-defined clusters in the MNIST plot indicate that the majority of the examples are separable via a kNN classifier in pixel space. The MNIST-1D plot, meanwhile, reveals a lack of well-defined clusters which suggests that learning a nonlinear 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)



Network

- 40 inputs
- 10 outputs
- Two hidden layers
 - 100 hidden units each
- SGD with batch size 100, learning rate 0.1
- 6000 steps (?? Epochs)

```
# 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)

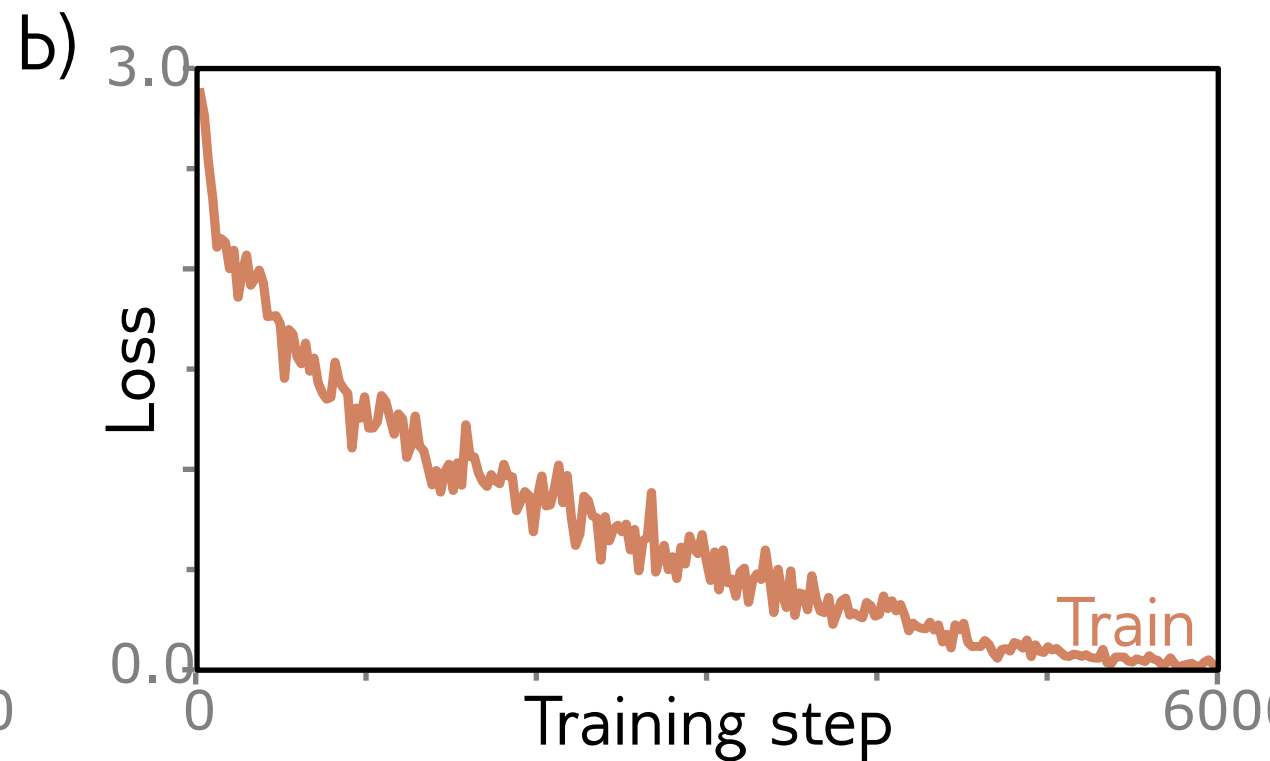
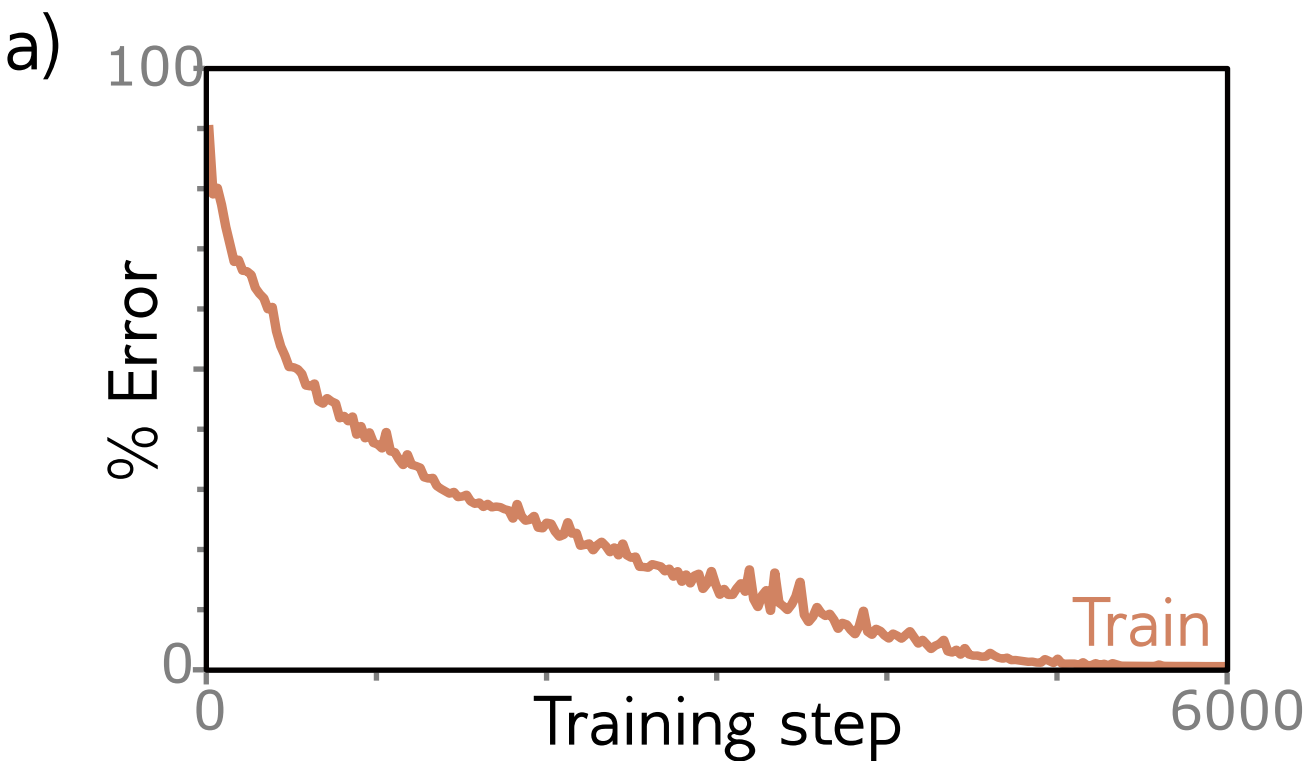
# object that decreases learning rate by half every 10 epochs
scheduler = StepLR(optimizer, step_size=10, gamma=0.5)

# load the data into a class that creates the batches
data_loader = DataLoader(TensorDataset(x_train,y_train), batch_size=100, shuffle=True)
```

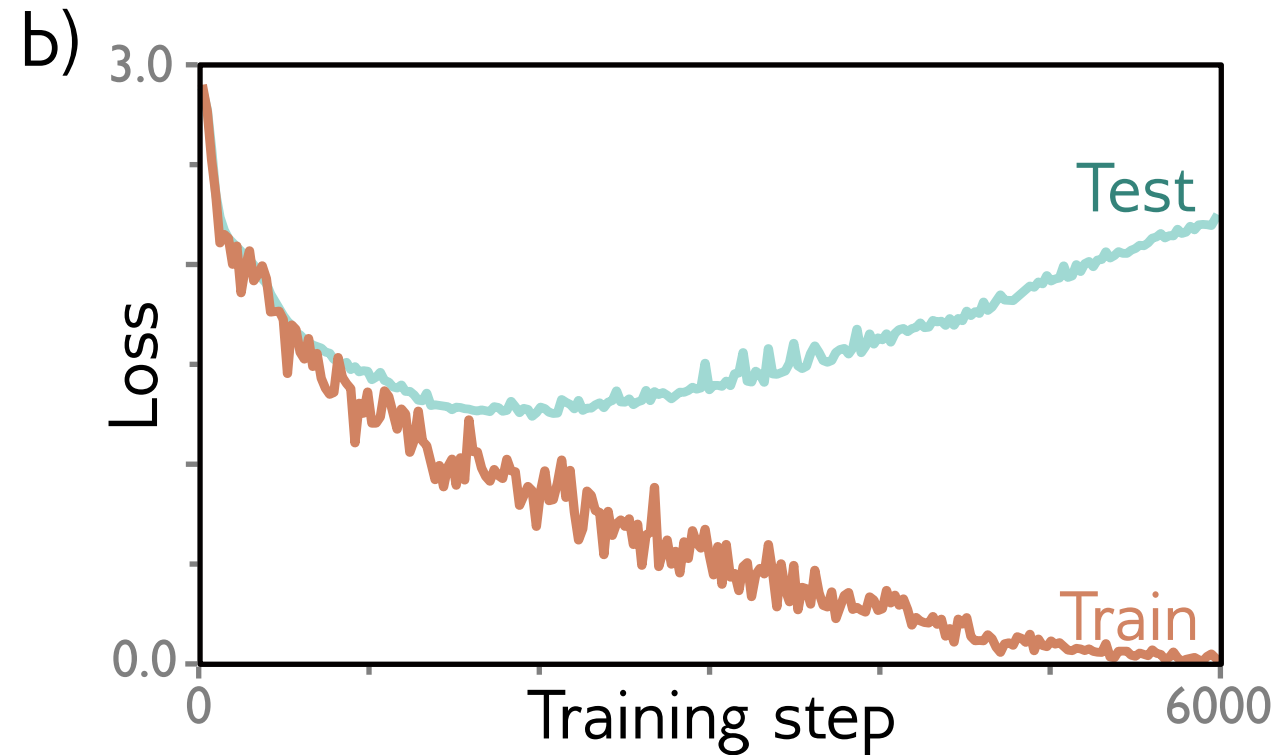
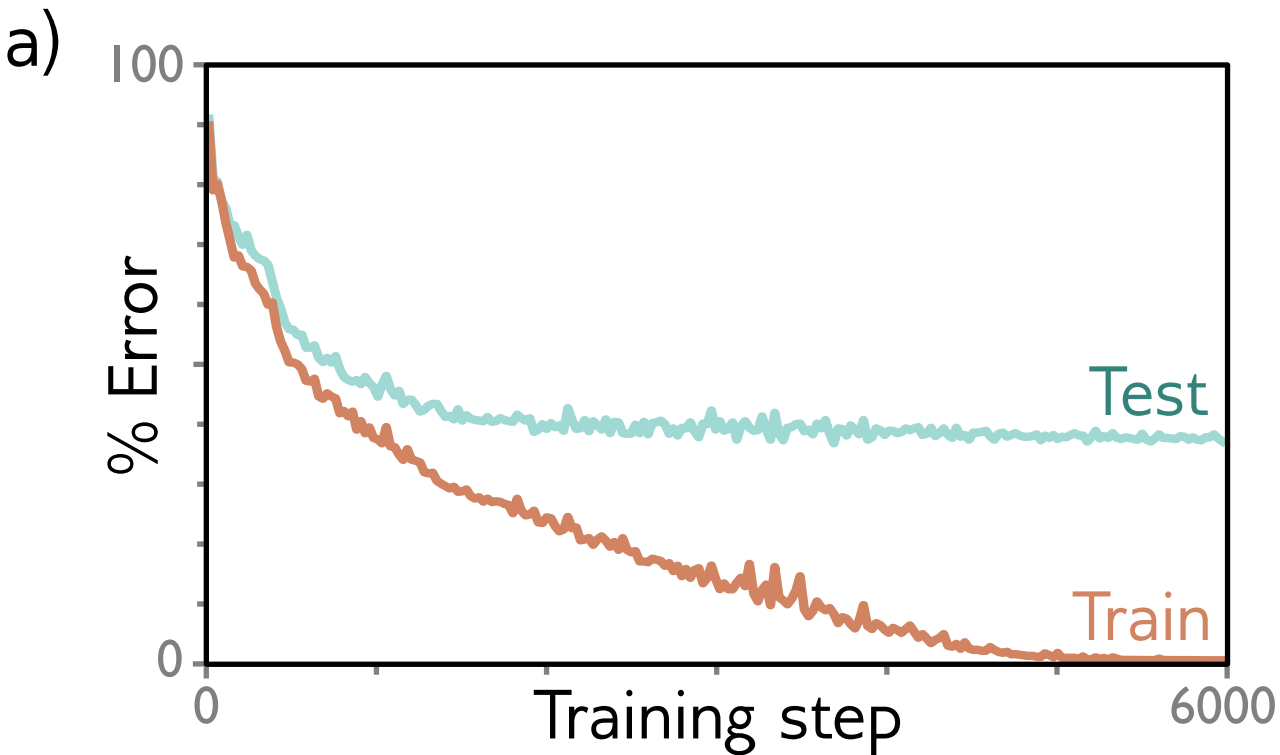
```
model = torch.nn.Sequential(
    torch.nn.Linear(40, 100),
    torch.nn.ReLU(),
    torch.nn.Linear(100, 100),
    torch.nn.ReLU(),
    torch.nn.Linear(100, 10))
```

```
=====
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



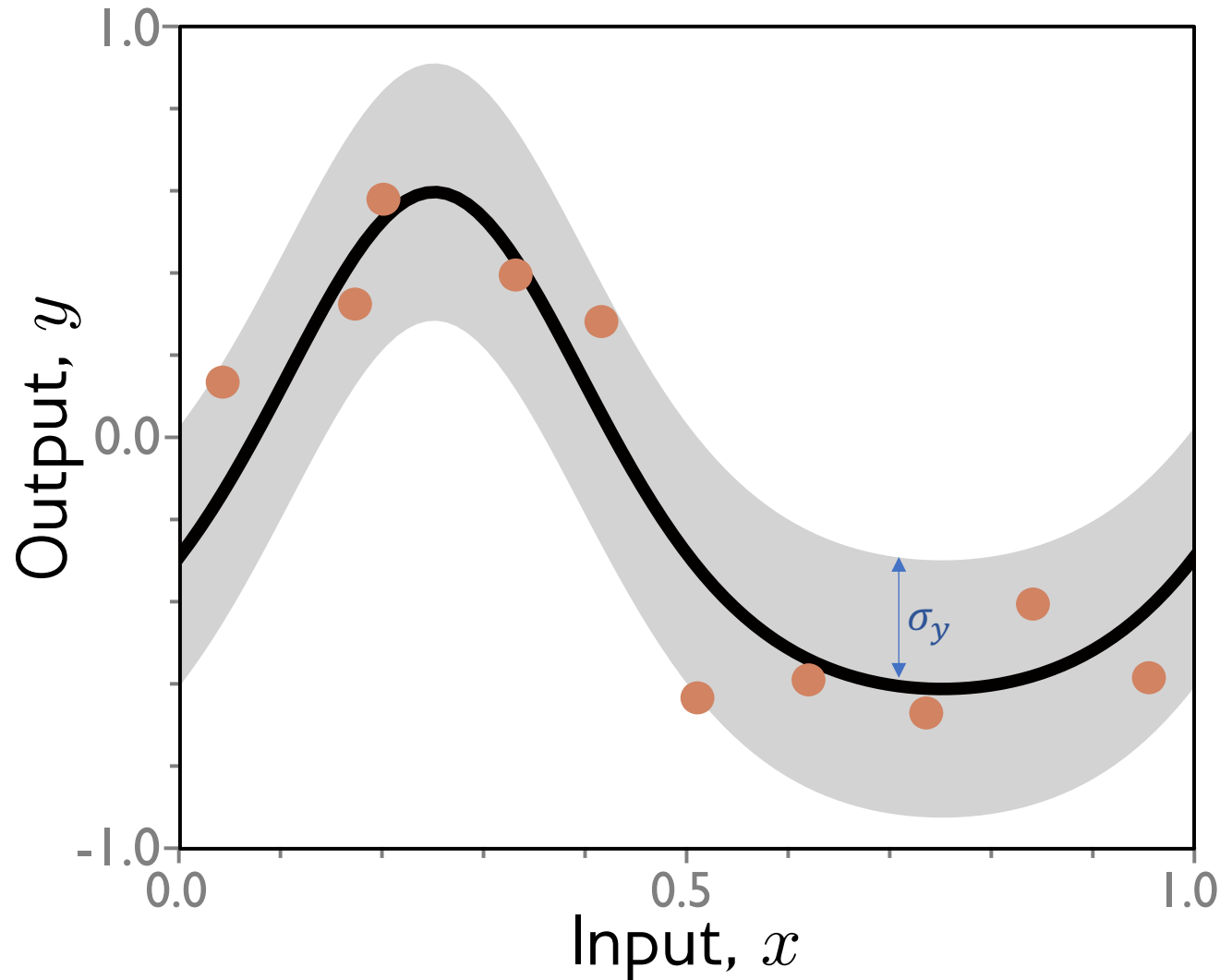
The model has not **generalized** well to the new data

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 with Toy Model

Dataset



"True" function:

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

Add small uniform noise to x :

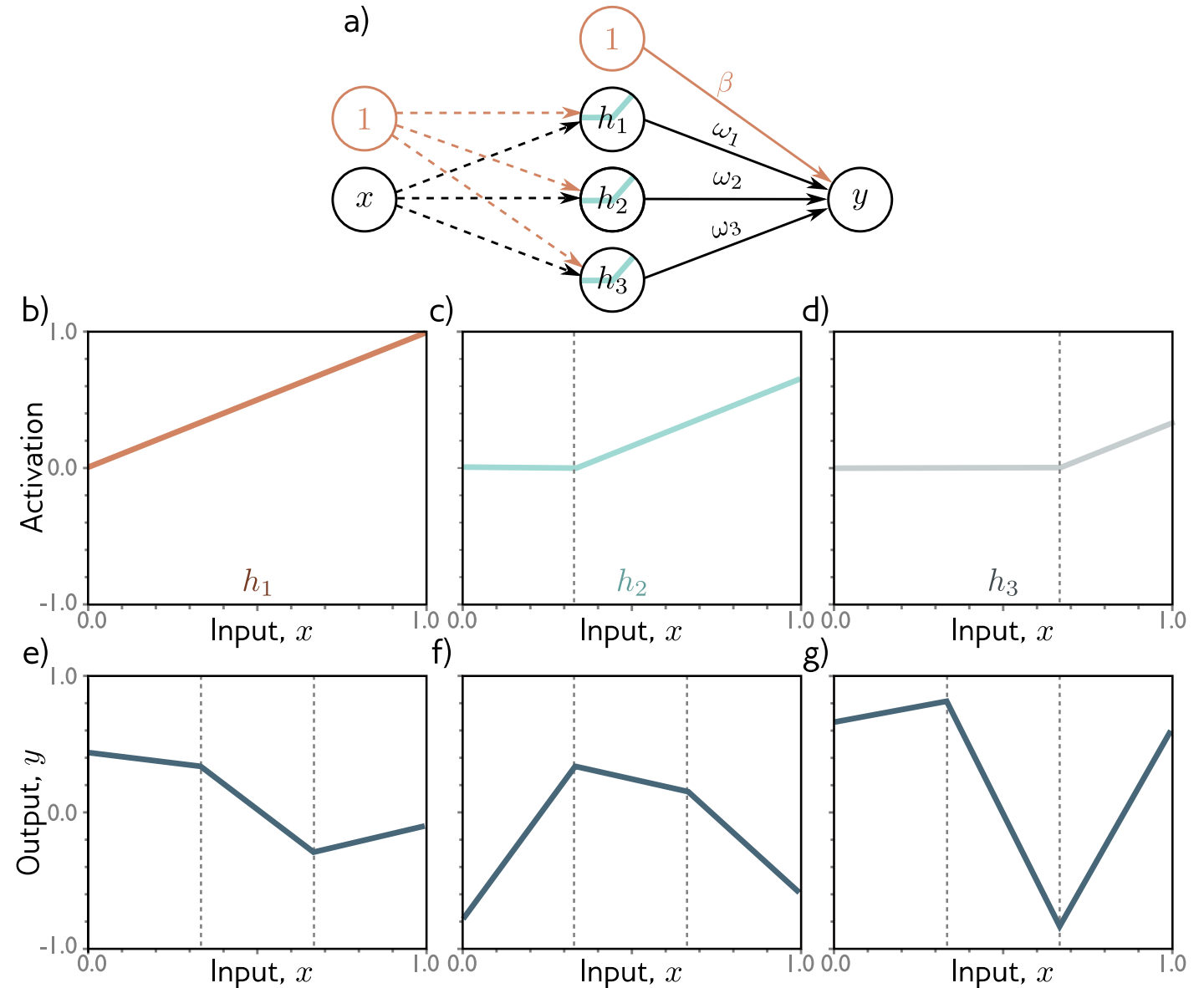
$$x = x + \mathcal{U}(\pm 1/\text{num_data})$$

Add small Gaussian noise to y :

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

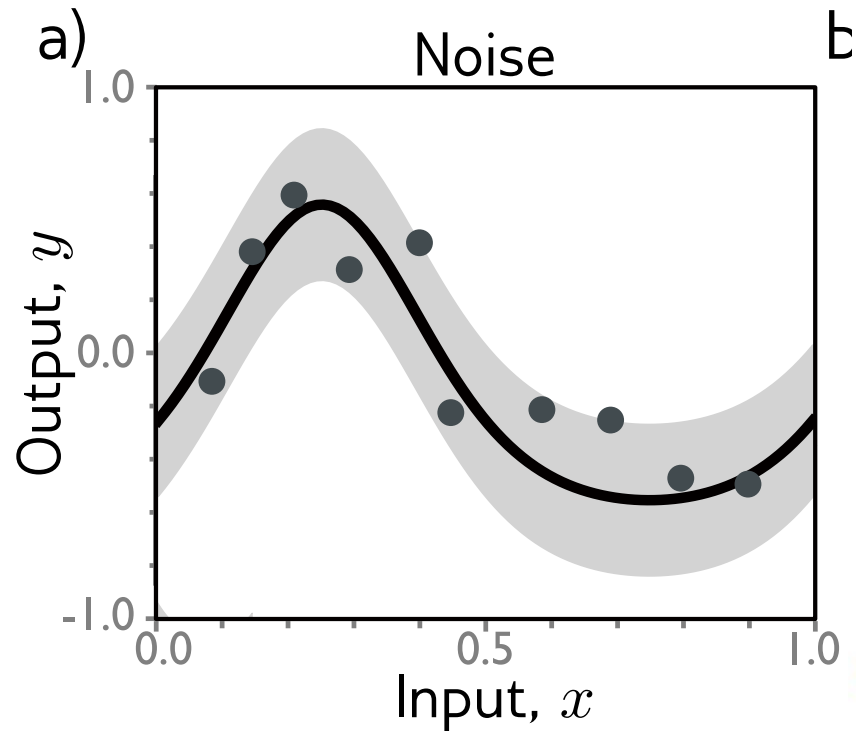
Toy model

- 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 soln in closed-form
- A piecewise linear model with D regions.

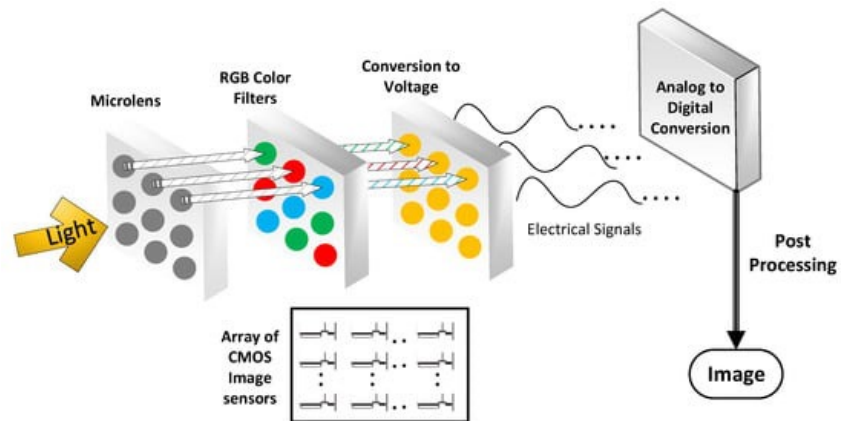


Three possible sources of error:
noise, bias and *variance*

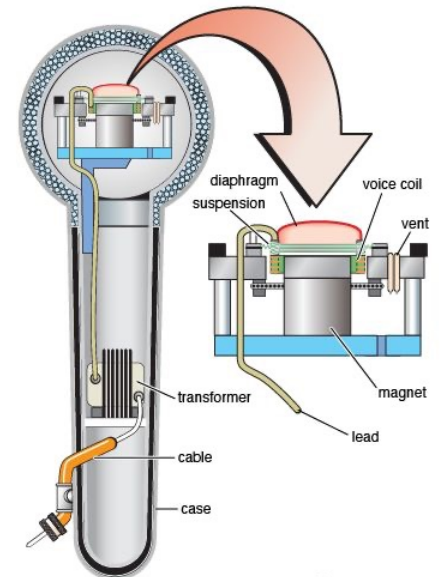
Noise, bias, and variance



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



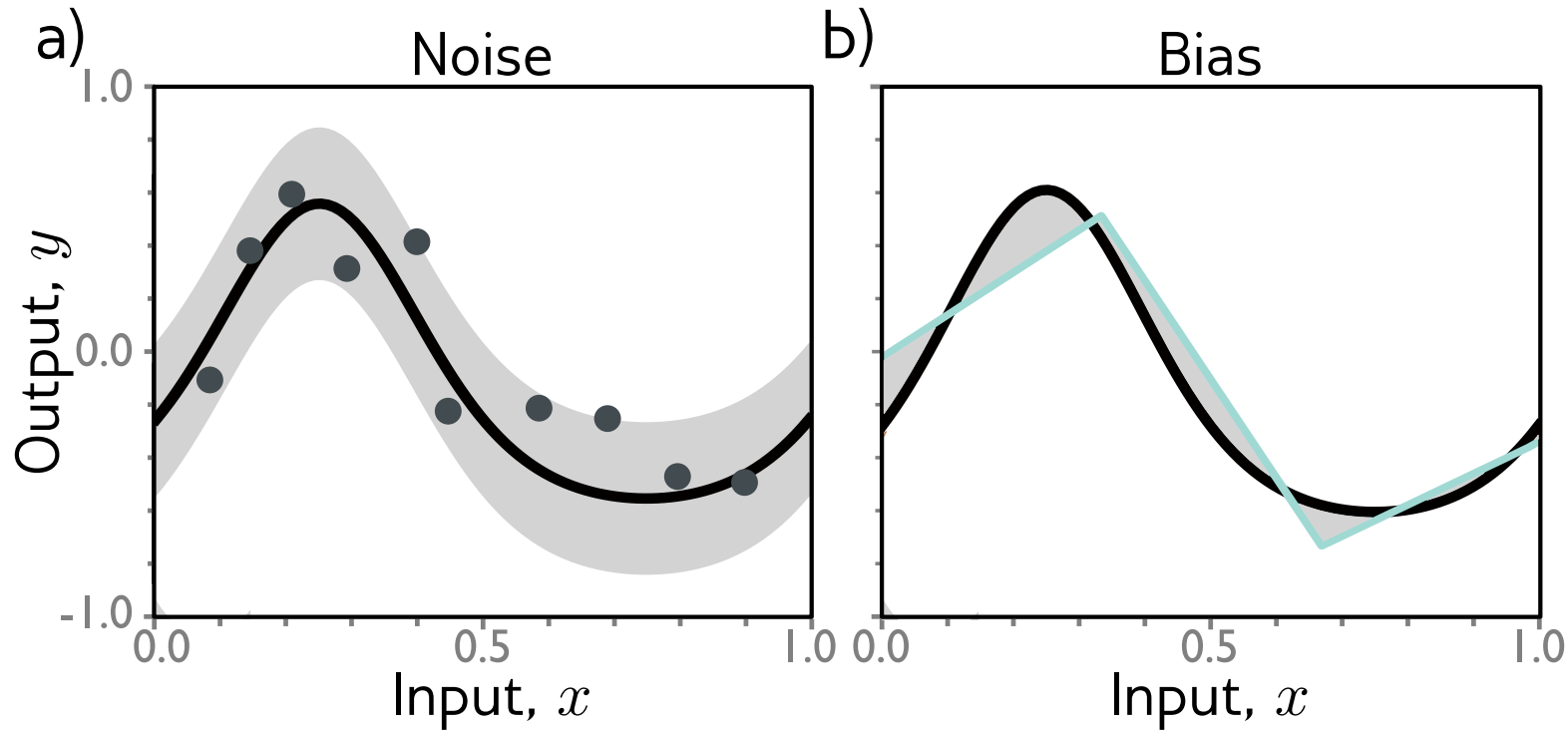
<https://images.app.goo.gl/2PuBhaFpfdL9Pyjb8>



© Merriam-Webster Inc.

<https://images.app.goo.gl/CMDaXSCdX4pgN8Yx7>

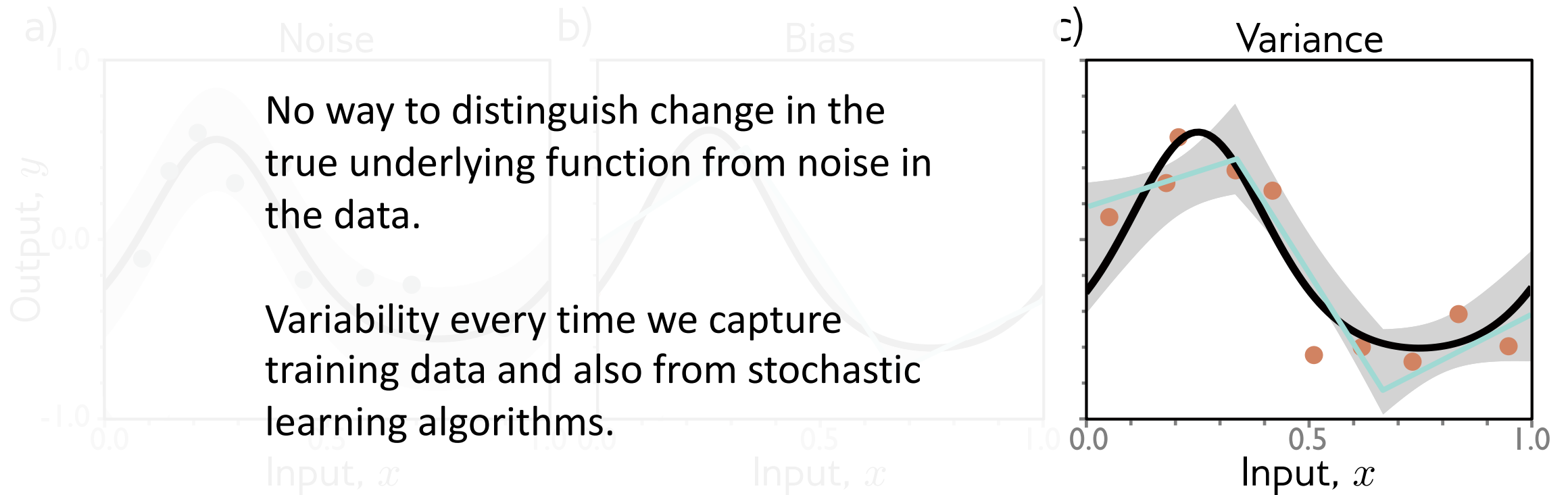
Noise, **bias**, and variance



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



Least squares regression only

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

- We can show that:

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

Expectation over noise
in training data

Expectation over
noise in test data

Actual model

Best possible model if
we had infinite data

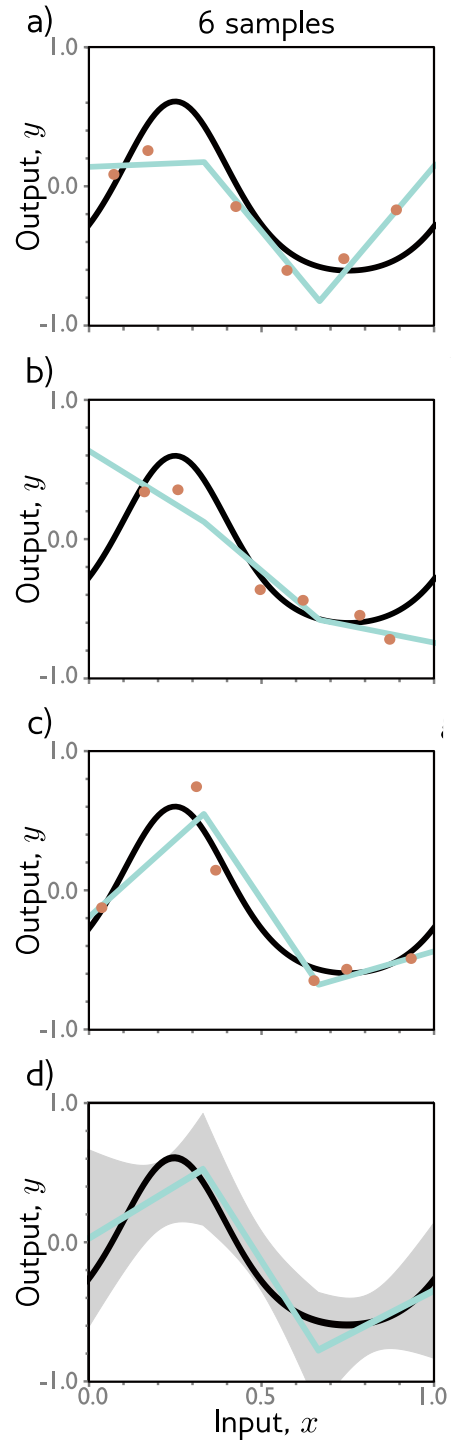
True function

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

Measuring performance

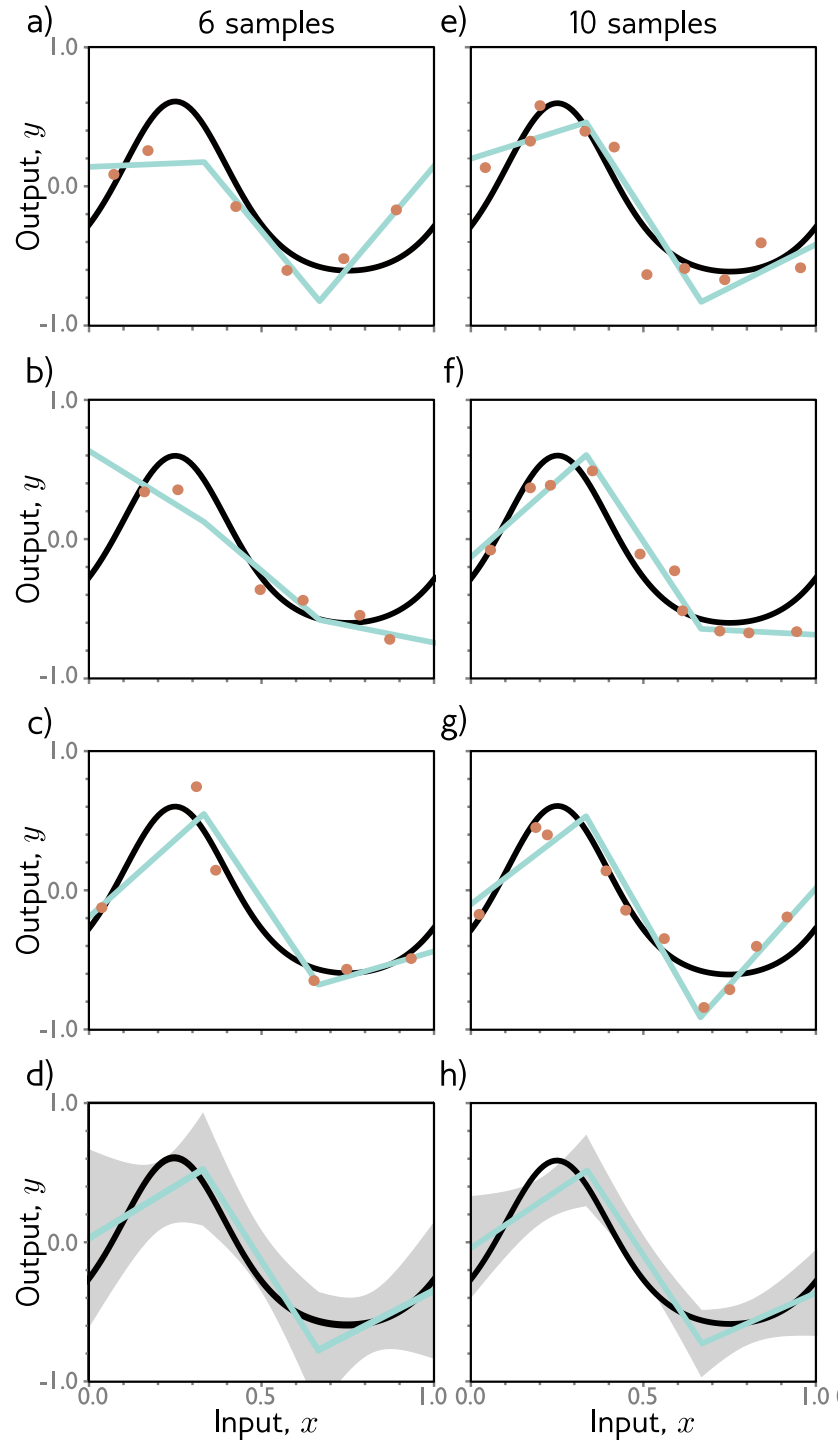
- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Curse of dimensionality & weird properties of high dimensional space
- Choosing hyperparameters

Variance



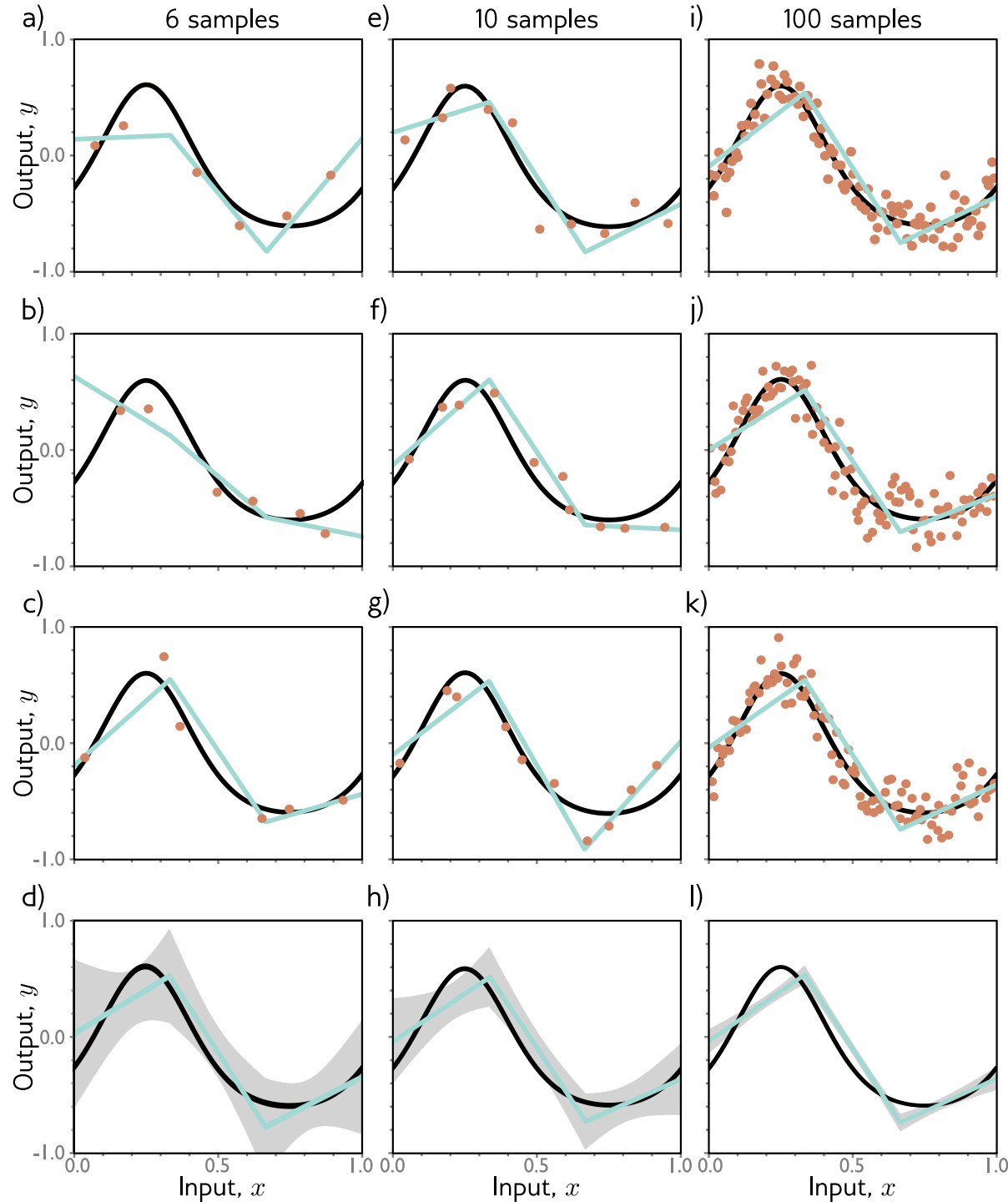
When measuring (capturing) 6 data samples and a fixed model (e.g. 3 hidden units), we get different optimal fits every time.

Variance



Can reduce
variance by
adding more
samples

Variance

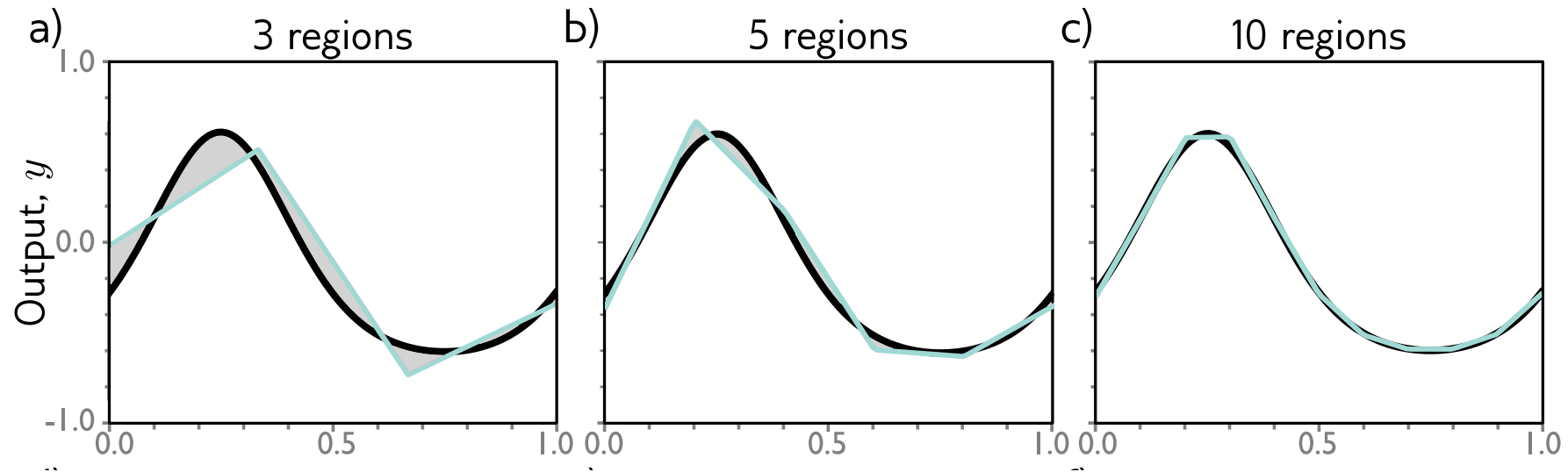


Can reduce
variance by
adding more
samples

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

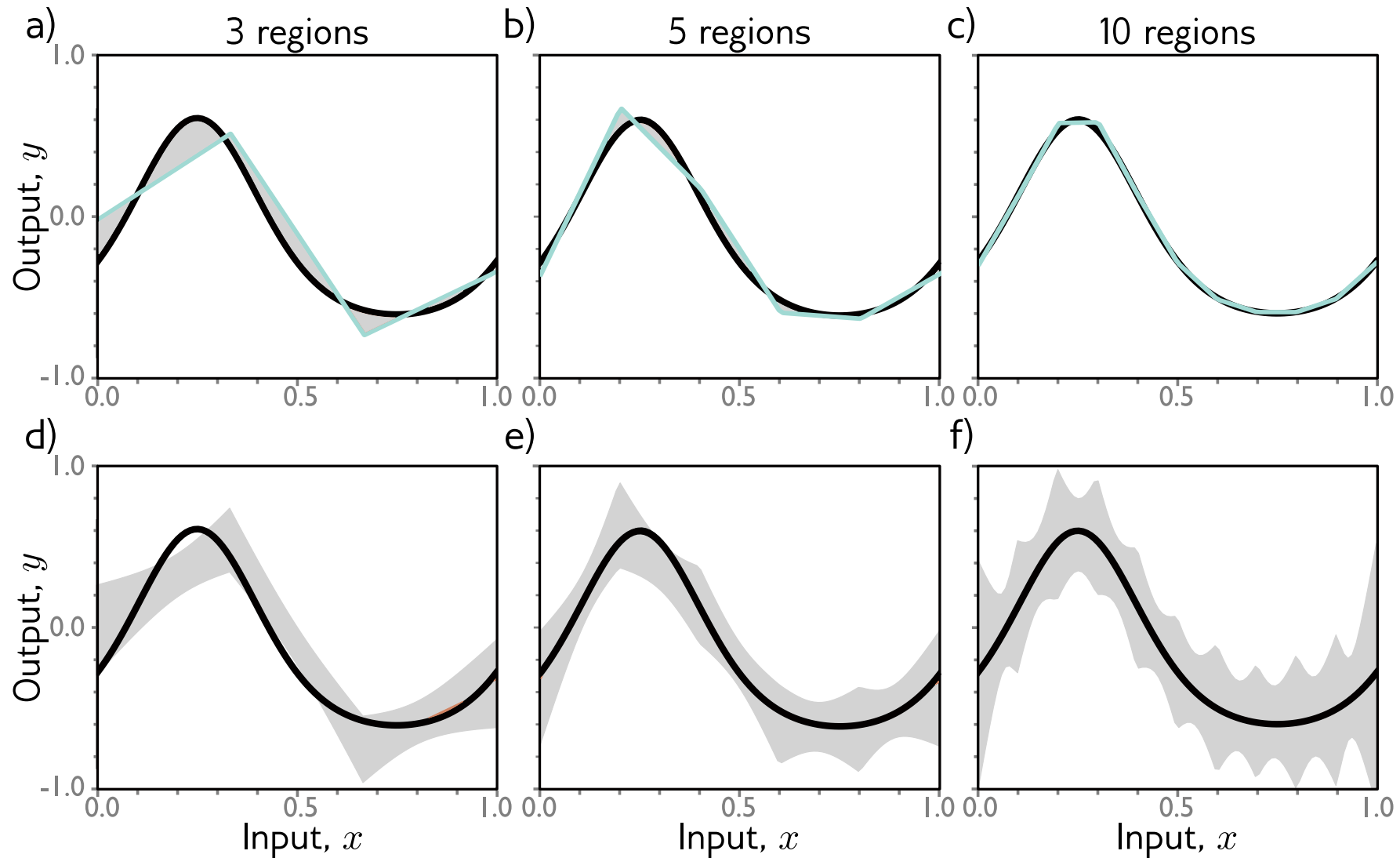
Reducing bias



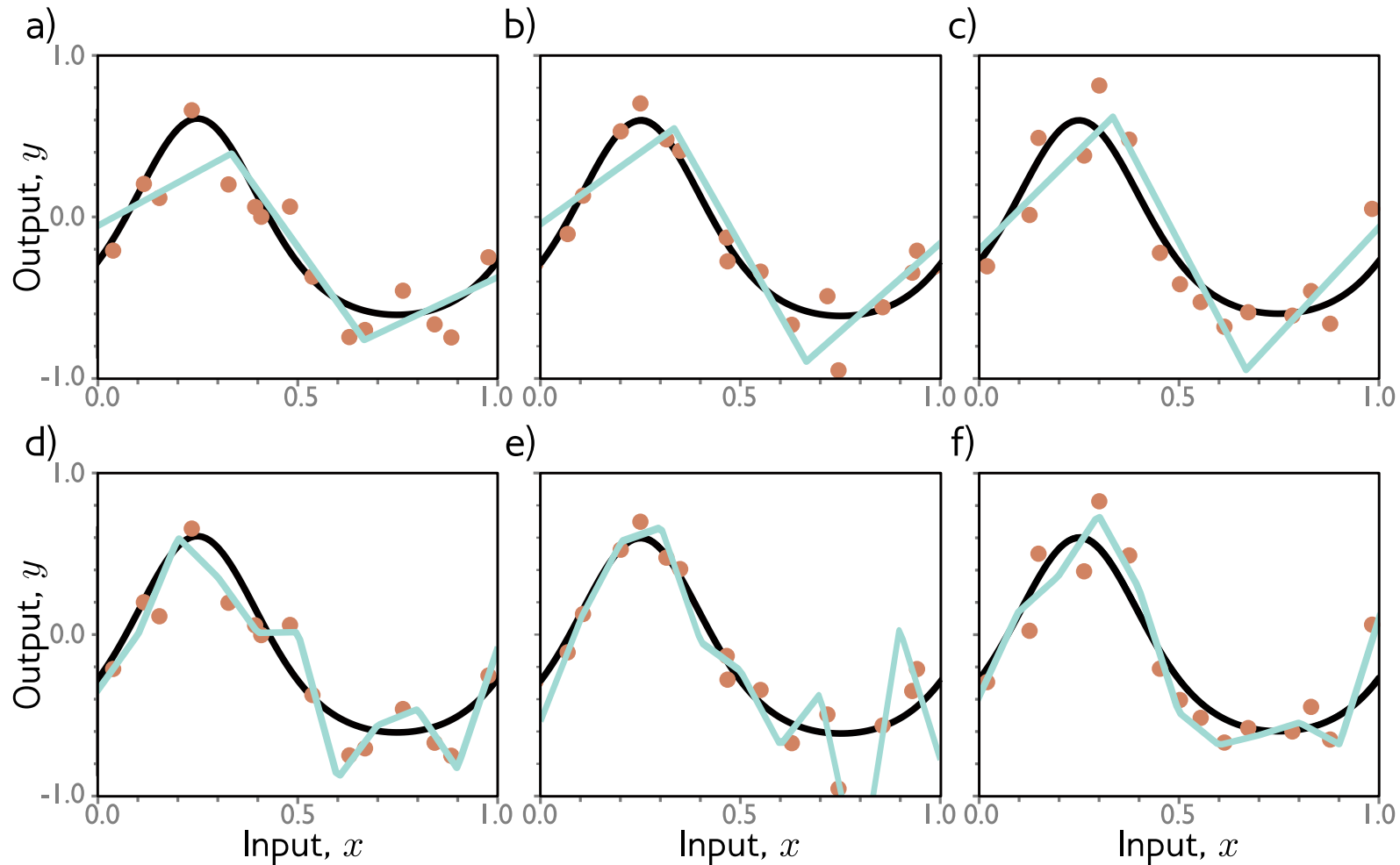
We can reduce bias by adding more model capacity.

In this case, adding more hidden units.

Reducing bias \rightarrow Increases variance!!

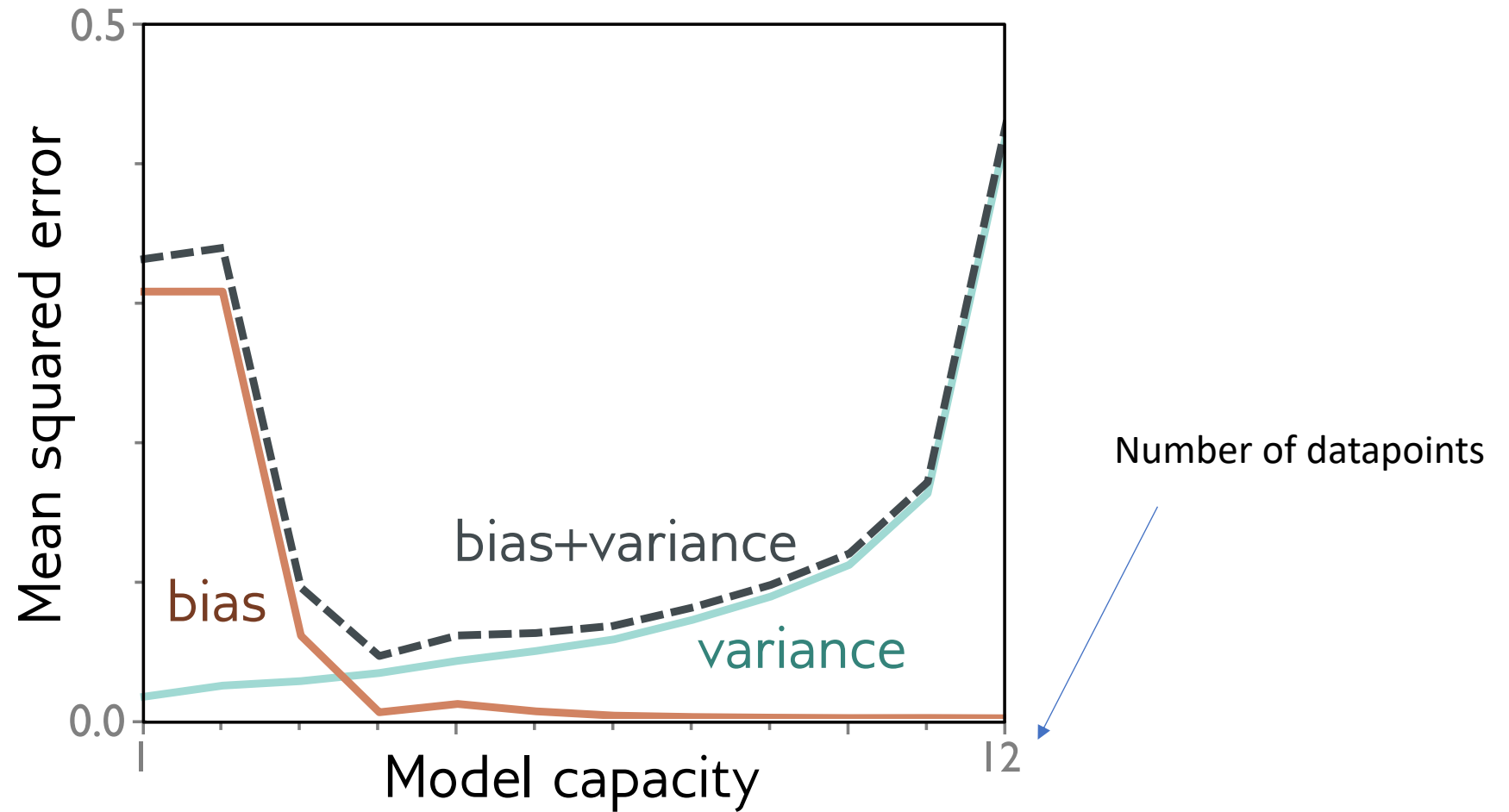


Why does variance increase? Overfitting



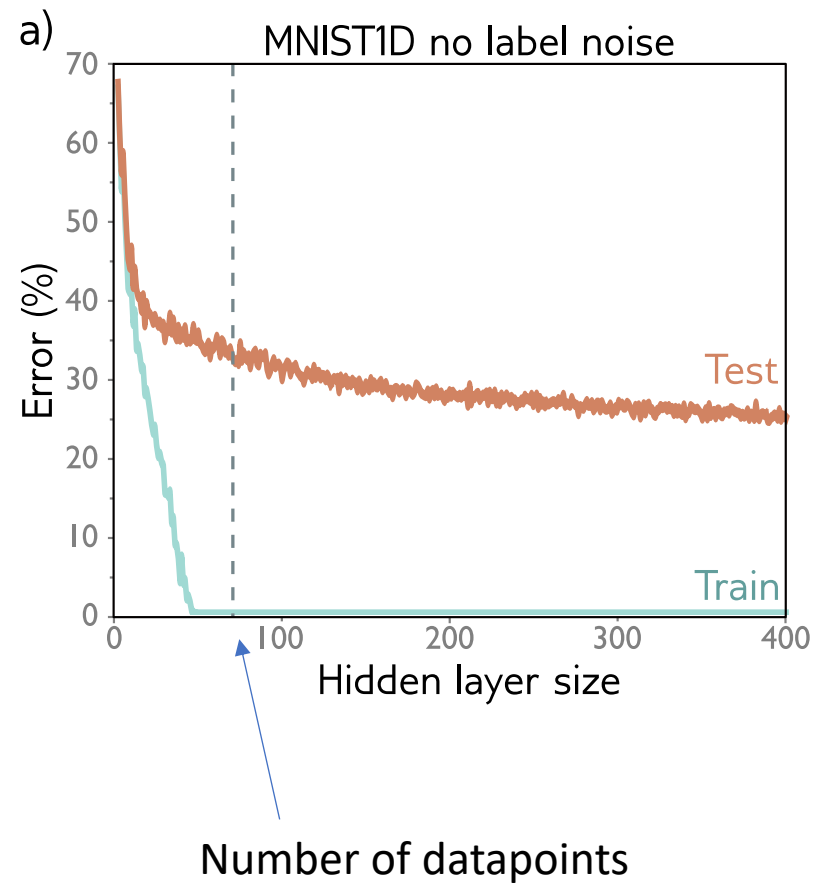
Describes the training data better, but not the true underlying function (black curve)

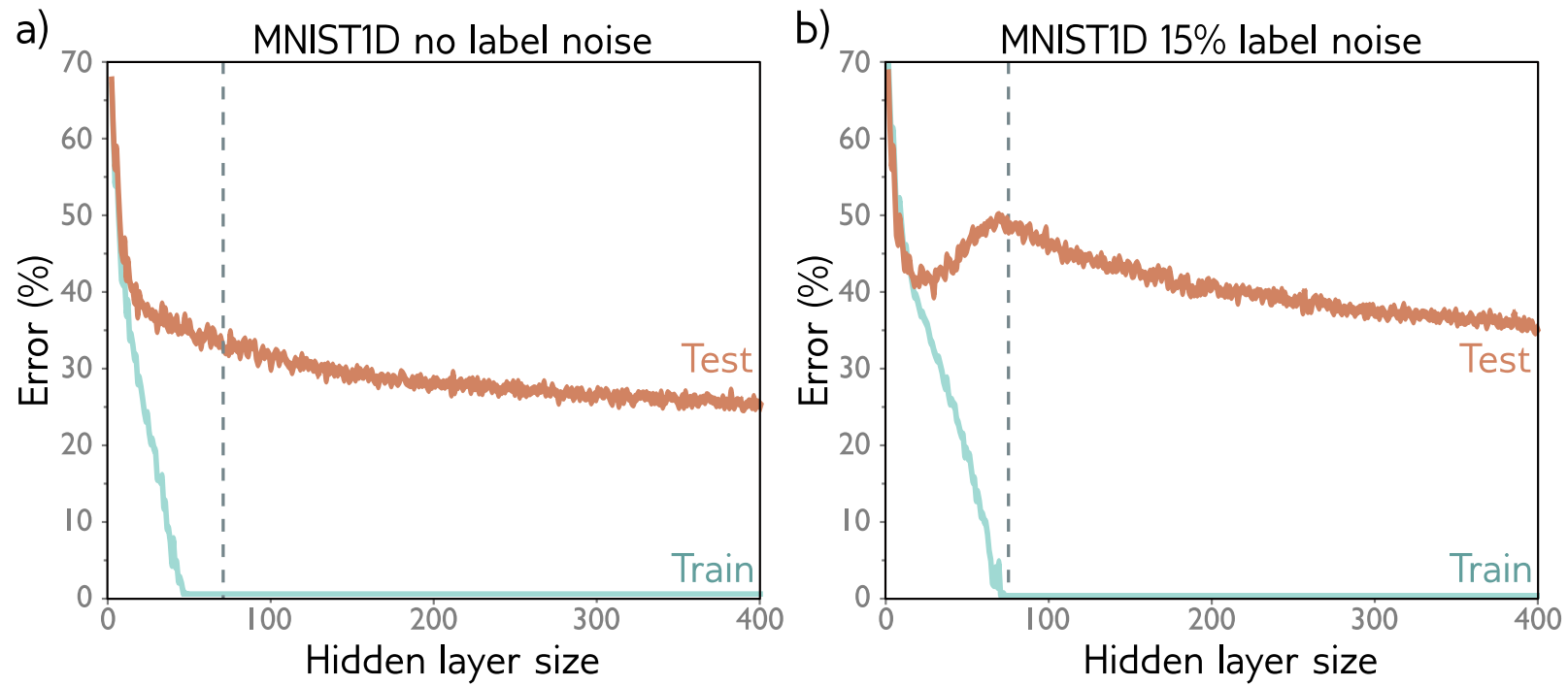
Bias and variance trade-off



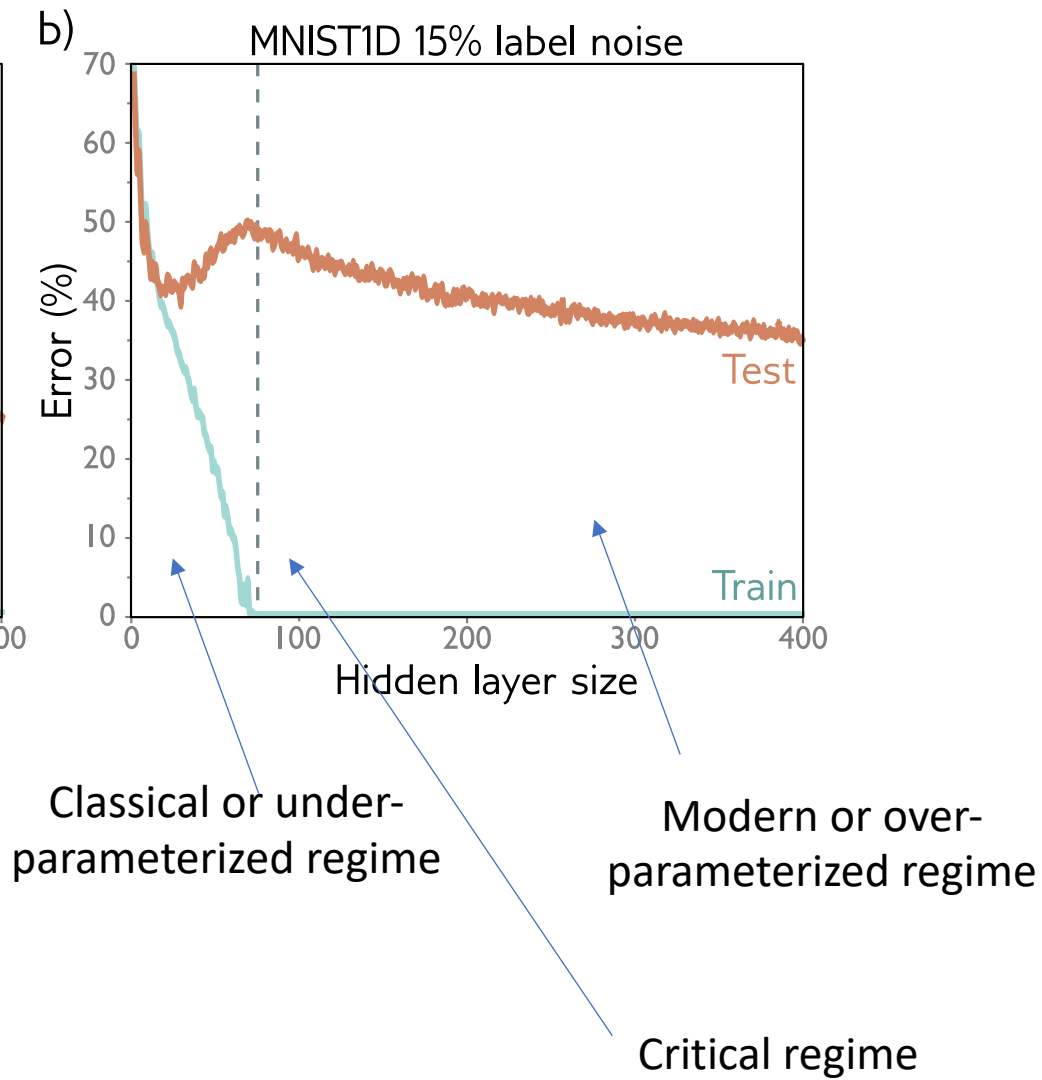
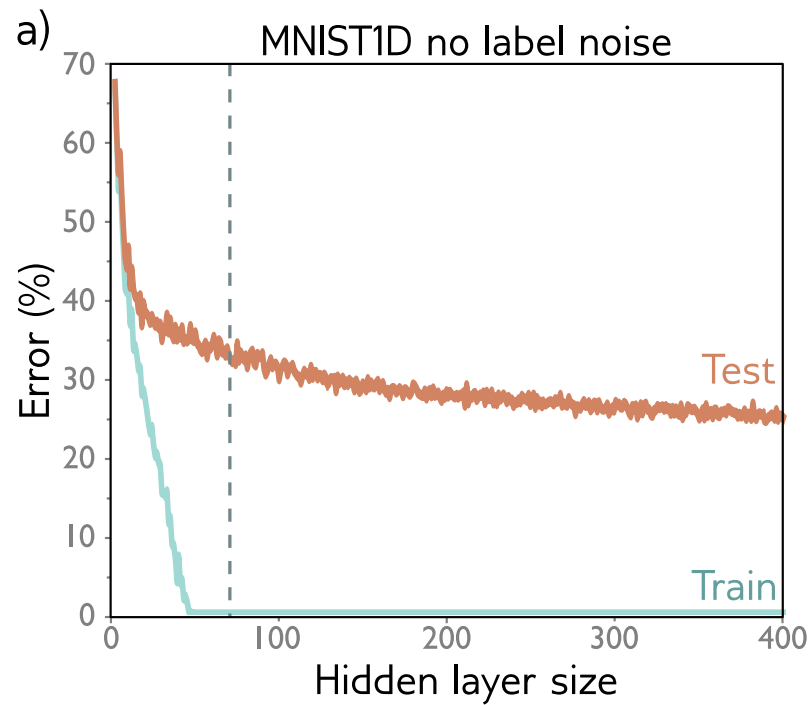
Measuring performance

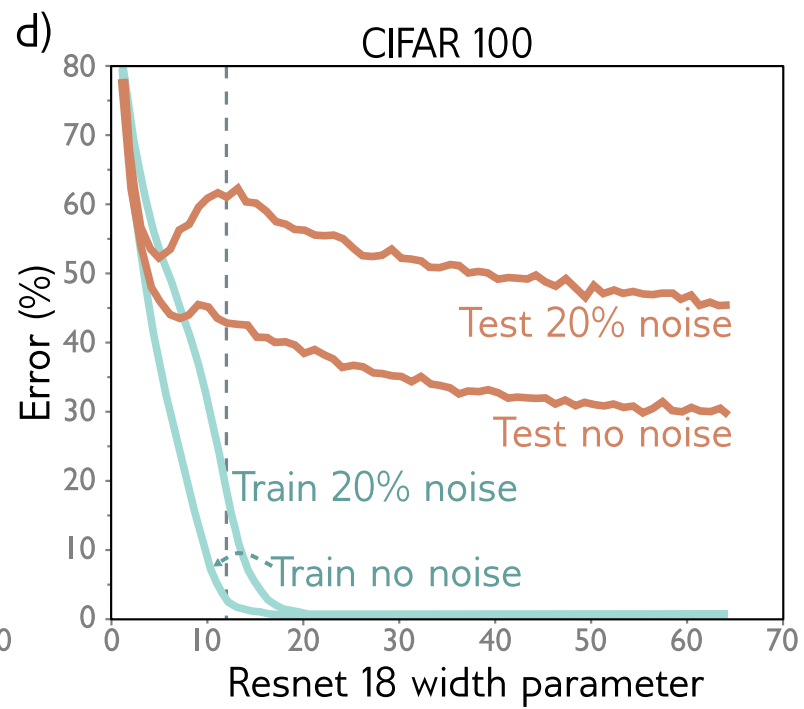
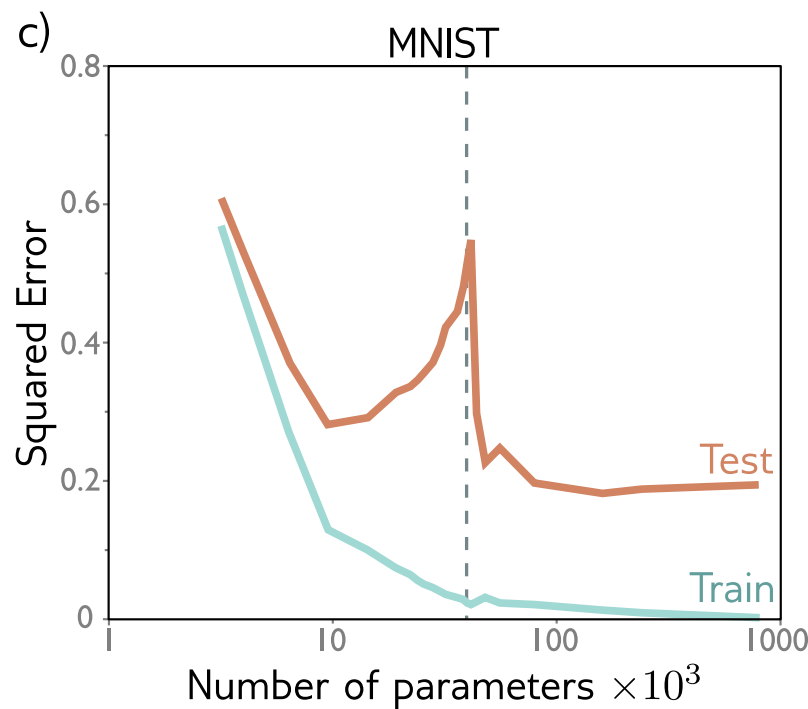
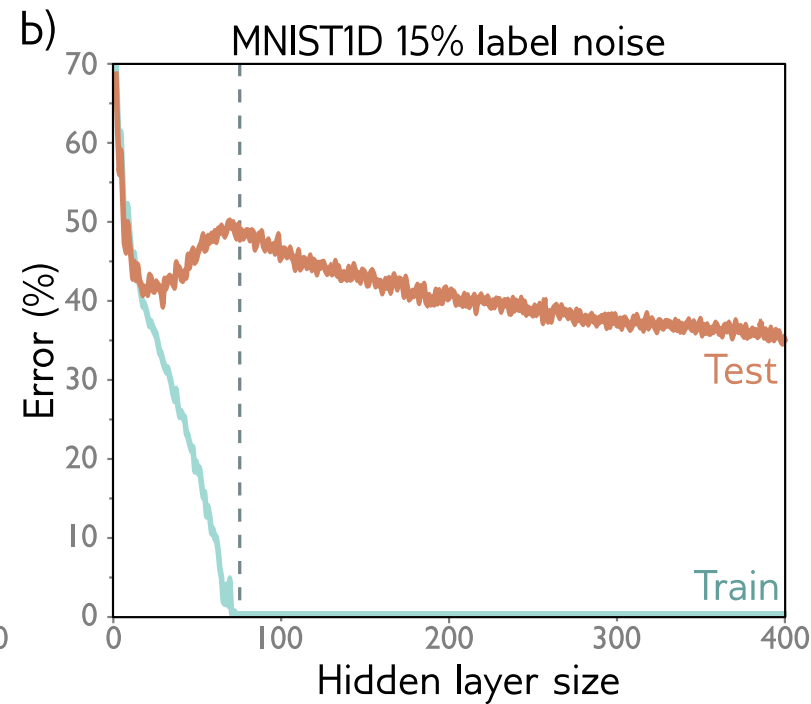
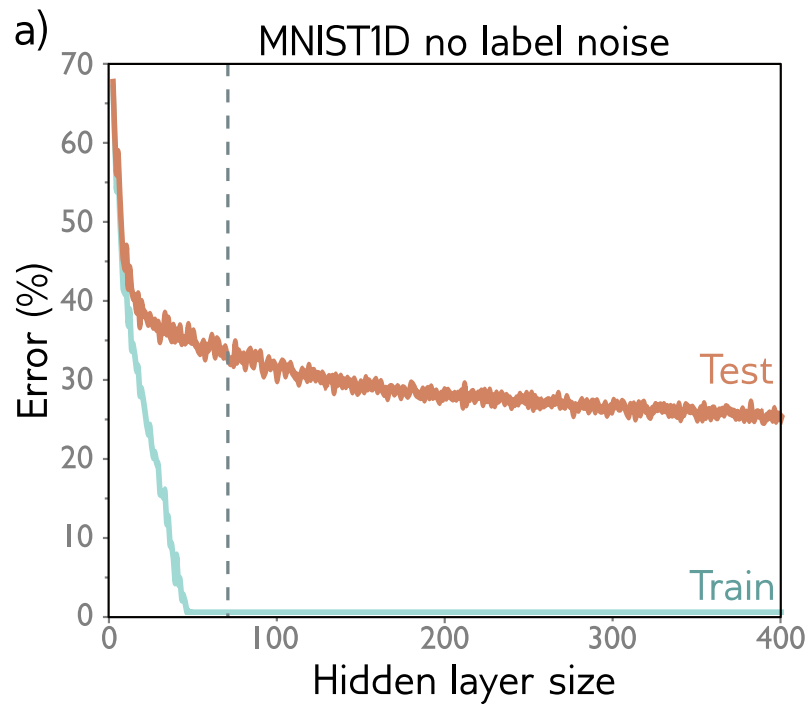
- MNIST1D dataset model and performance
- Noise, bias, and variance
- Reducing variance
- Reducing bias & bias-variance trade-off
- Double descent
- Curse of dimensionality & weird properties of high dimensional space
- Choosing hyperparameters

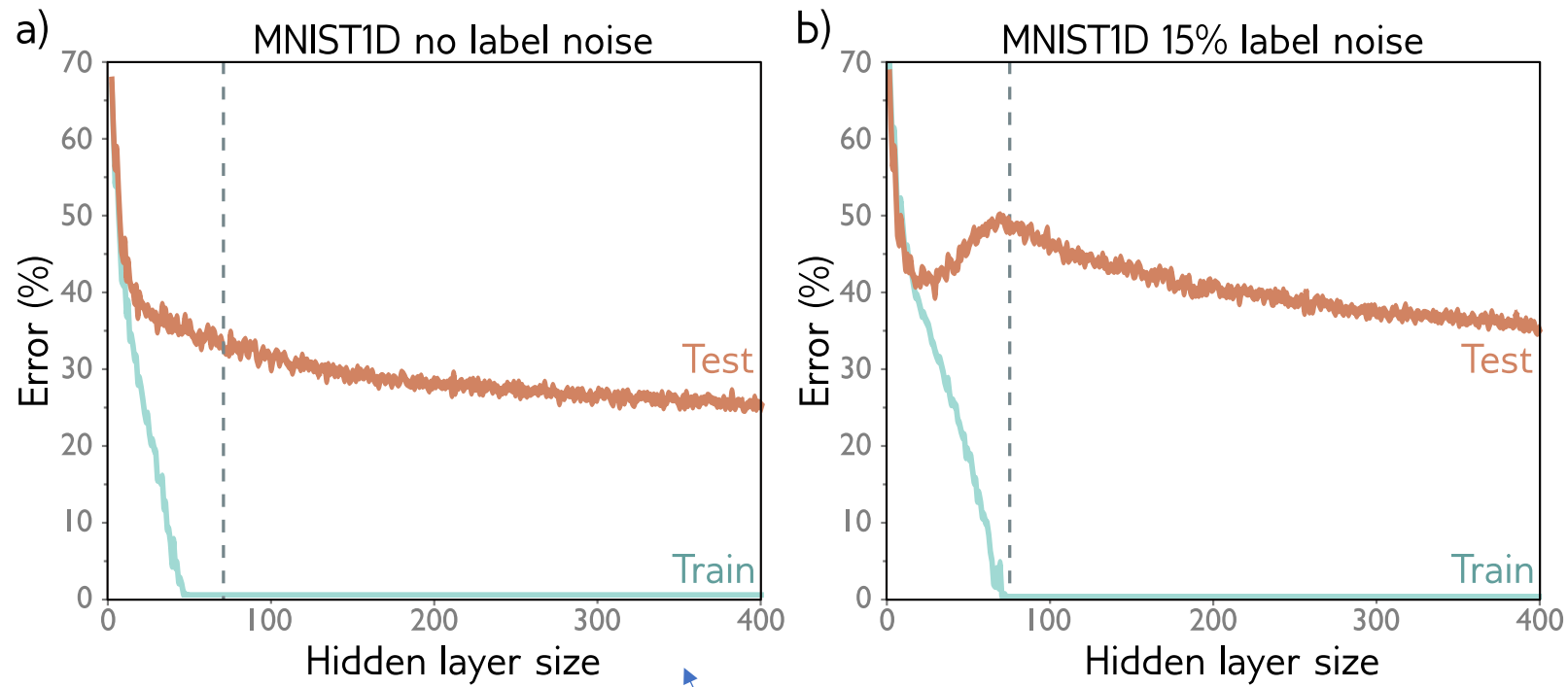




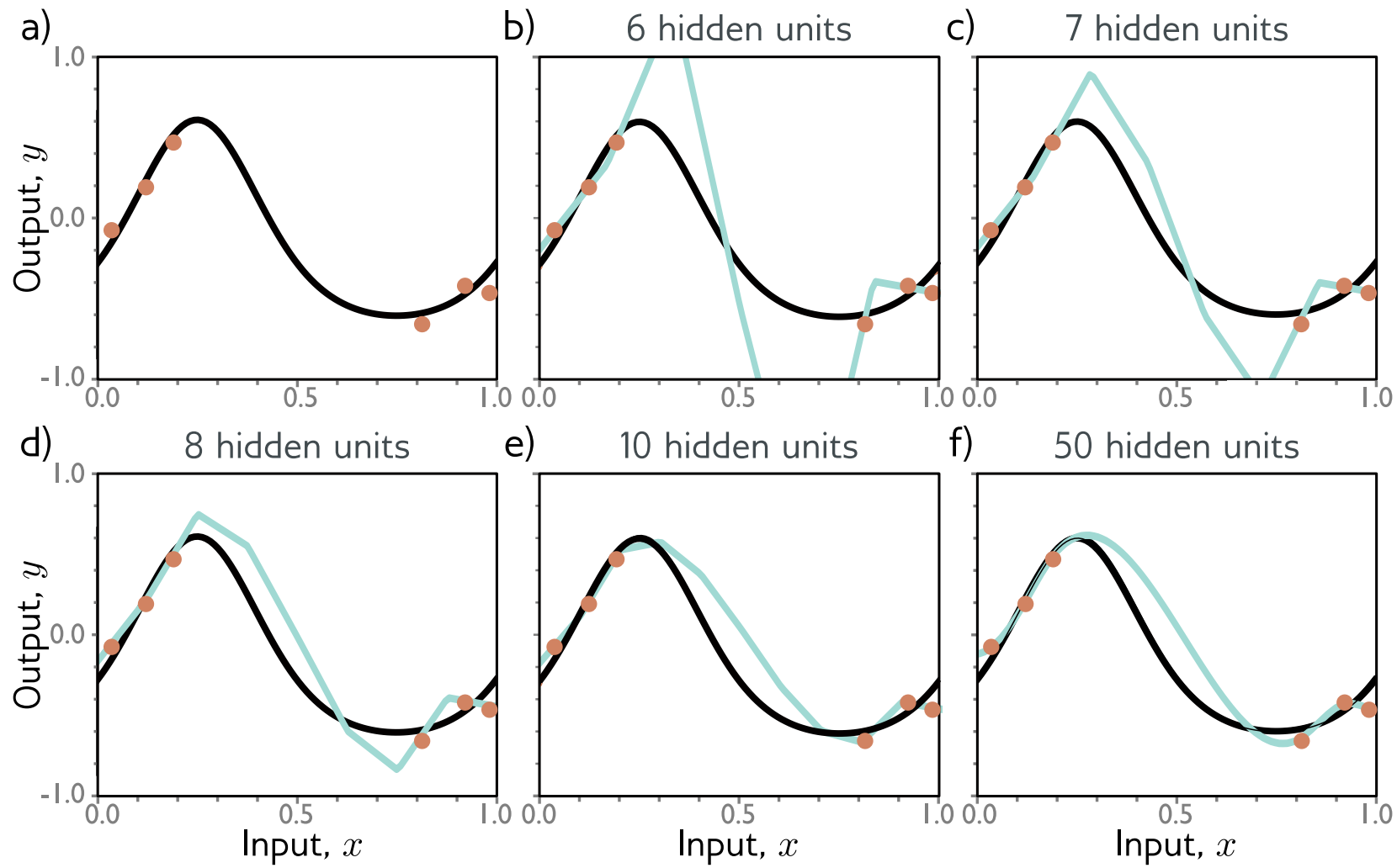
Double descent







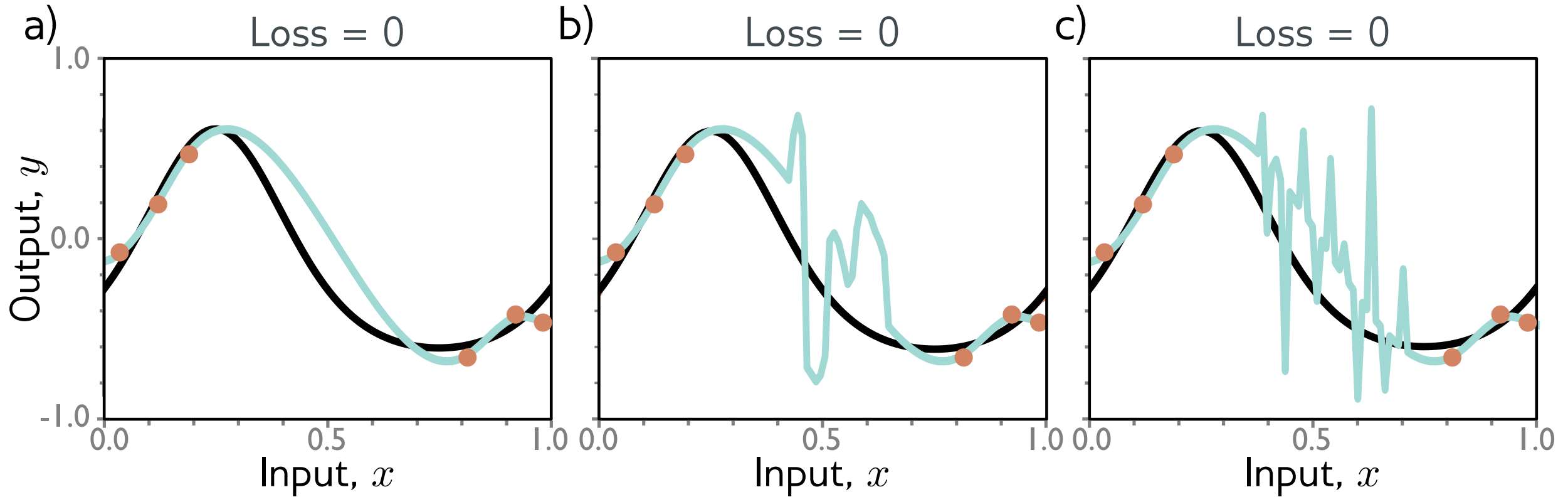
- Note that train data is very close to zero.
- Whatever is happening isn't happening at training data points
- Must be happening between the data points??



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**

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

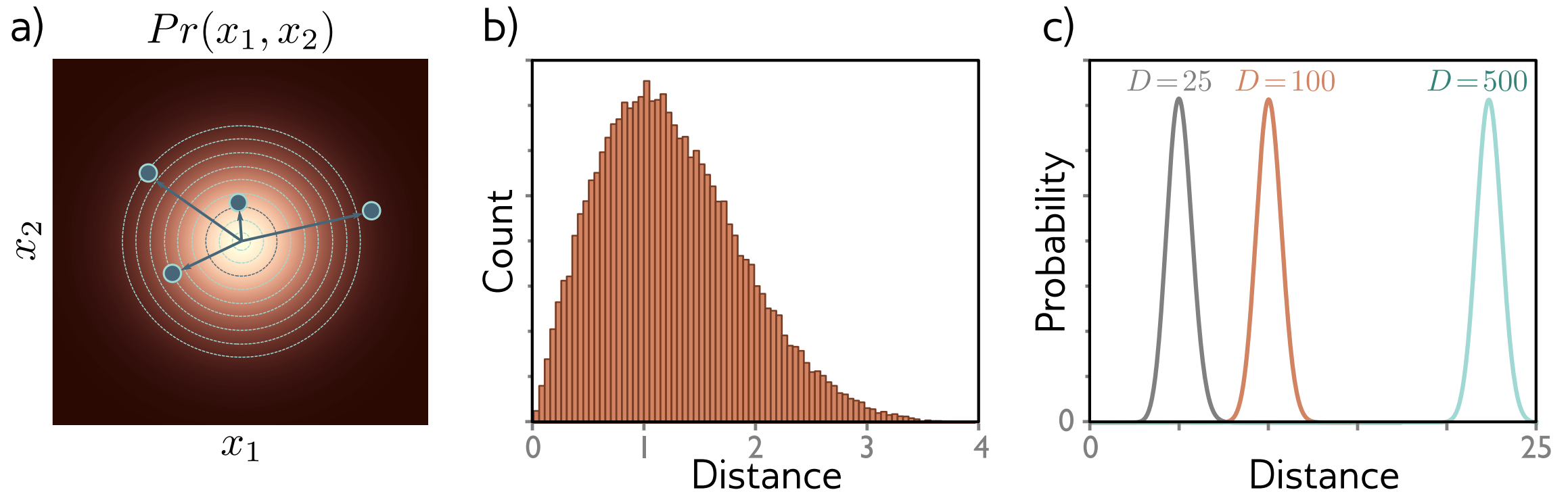
Curse of dimensionality

- 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**

Weird properties of high-dimensional space

- Two randomly sampled data points from normal are at right angles to each other with high likelihood
- Distance from the origin of random samples is roughly constant

Distance from the origin of random samples is roughly constant



Weird properties of high-dimensional space

- Two randomly sampled data points from normal are at right angles to each other with high likelihood
- Distance from the origin of random samples is roughly constant
- Most of the volume of a high dimensional orange is in the peel not in the pulp
- Volume of a diameter one hypersphere becomes zero
- Generate random points uniformly in hypercube, ratio of nearest to farthest becomes close to one.

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

Feedback?

