# Deep Learning for Data Science DS 542

https://dl4ds.github.io/fa2025/
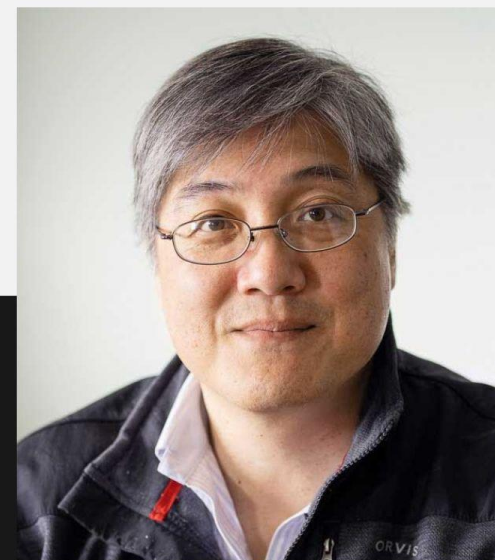
Shallow Neural Networks

# Announcements

Shared Compute Cluster (SCC) Tutorial next class (9/22)

- Bring your laptop next time!

- Will walk through account setup and ways to access the SCC.

Alumni Weekend Computer Science Distinguished Lecture

**Do LLMs Contain Concepts?**

with Prof. David Bau of Northeastern University

How do large language models think? Do they contain "concepts?" In this talk we will examine the internal mechanisms of LLM when performing several kinds of reasoning.
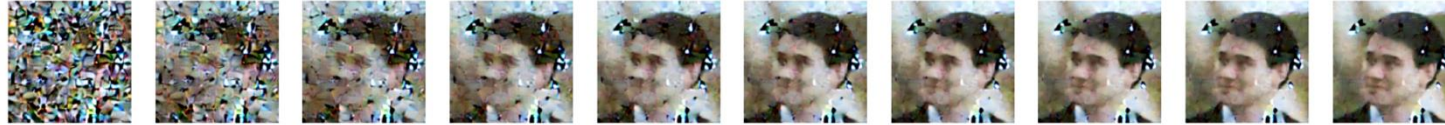
Sept 25th | 11am | CDS 1750

BOSTON UNIVERSITY

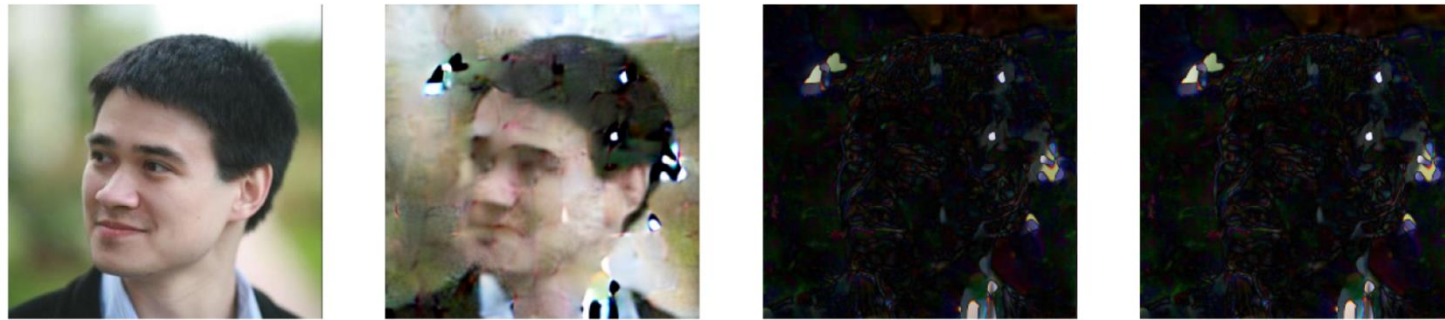**Boston University** College of Arts & Sciences Department of Computer Science

# Homework 2 FAQ

```python
# plot 10 of the images to show the progress
plot_images(decoded_images[num_images//10-1::(num_images+9)//10,:,:])
```
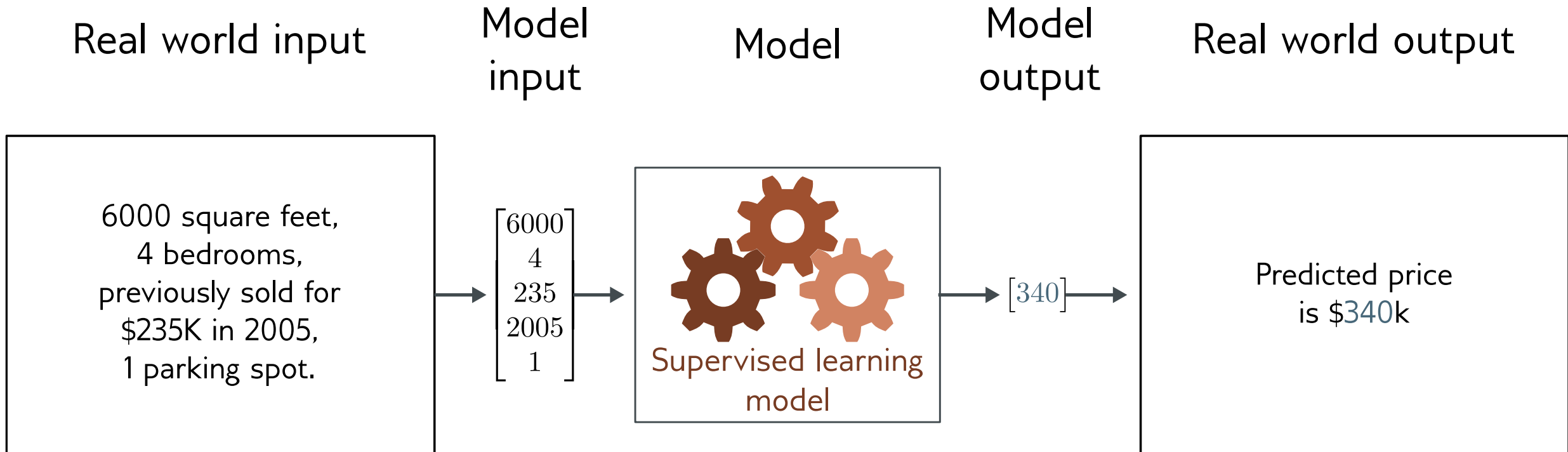


```python
last_image = decoded_image[-1:,:,:,:]

diff_image = (target_image - last_image).abs()

comparison_images = torch.cat([target_image, last_image, diff_image, diff_image / diff_image.max()], dim=0)
plot_images(comparison_images)
```
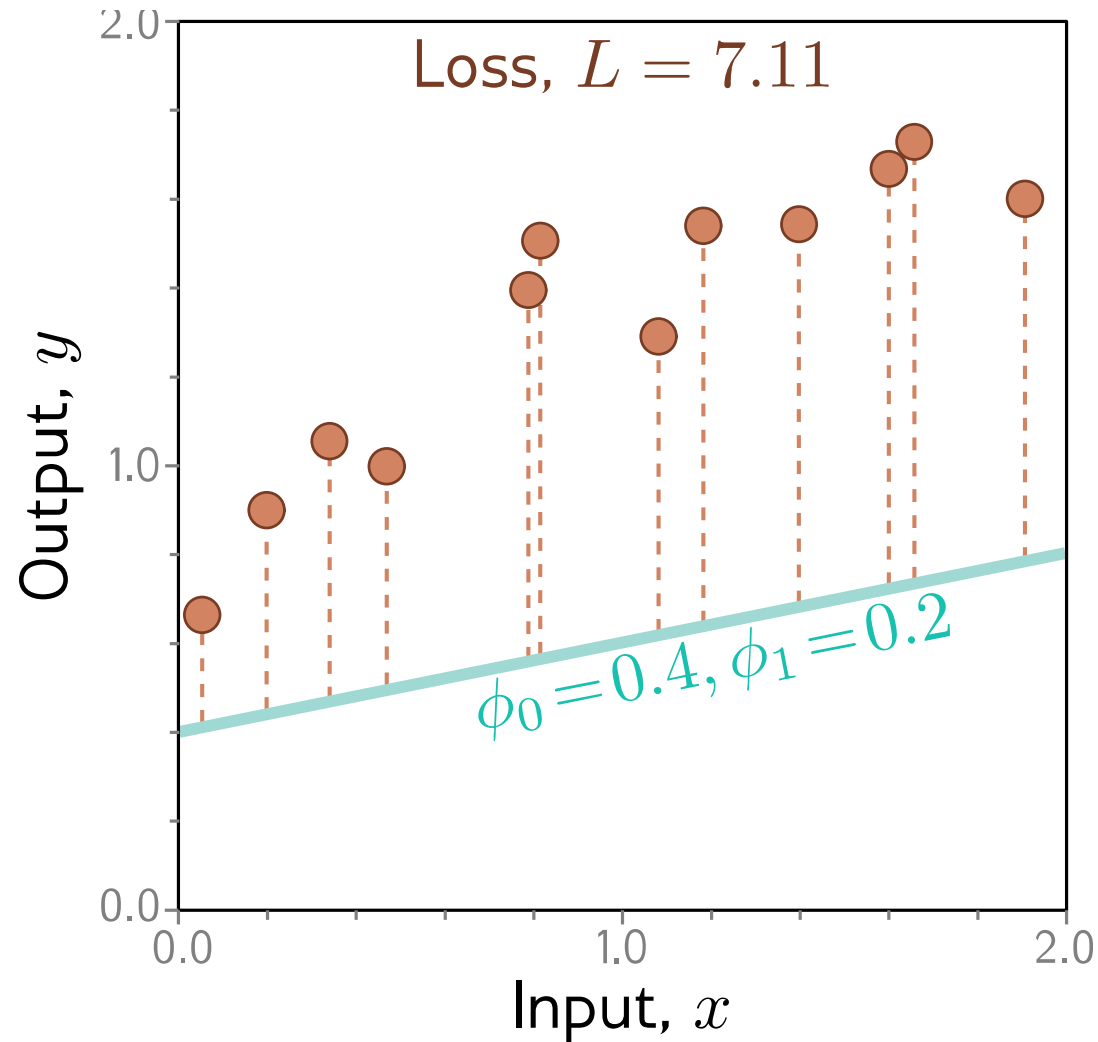
# Recap: Regression

**Real world input**     **Model input**     **Model**     **Model output**     **Real world output**

6000 square feet,
4 bedrooms,
previously sold for
$235K in 2005,
1 parking spot.

$$\begin{bmatrix} 6000 \\ 4 \\ 235 \\ 2005 \\ 1 \end{bmatrix}$$

Supervised learning model

$[340]$

Predicted price is $340k

- Univariate regression problem (one output, real value)
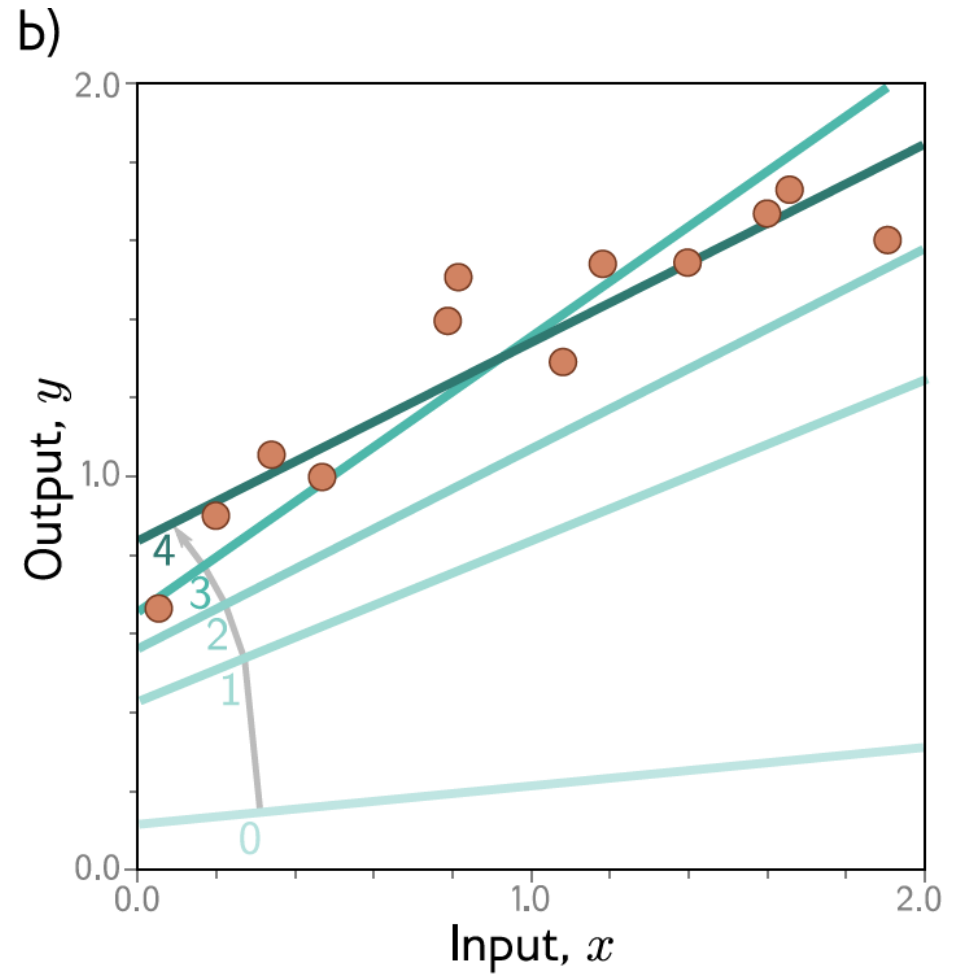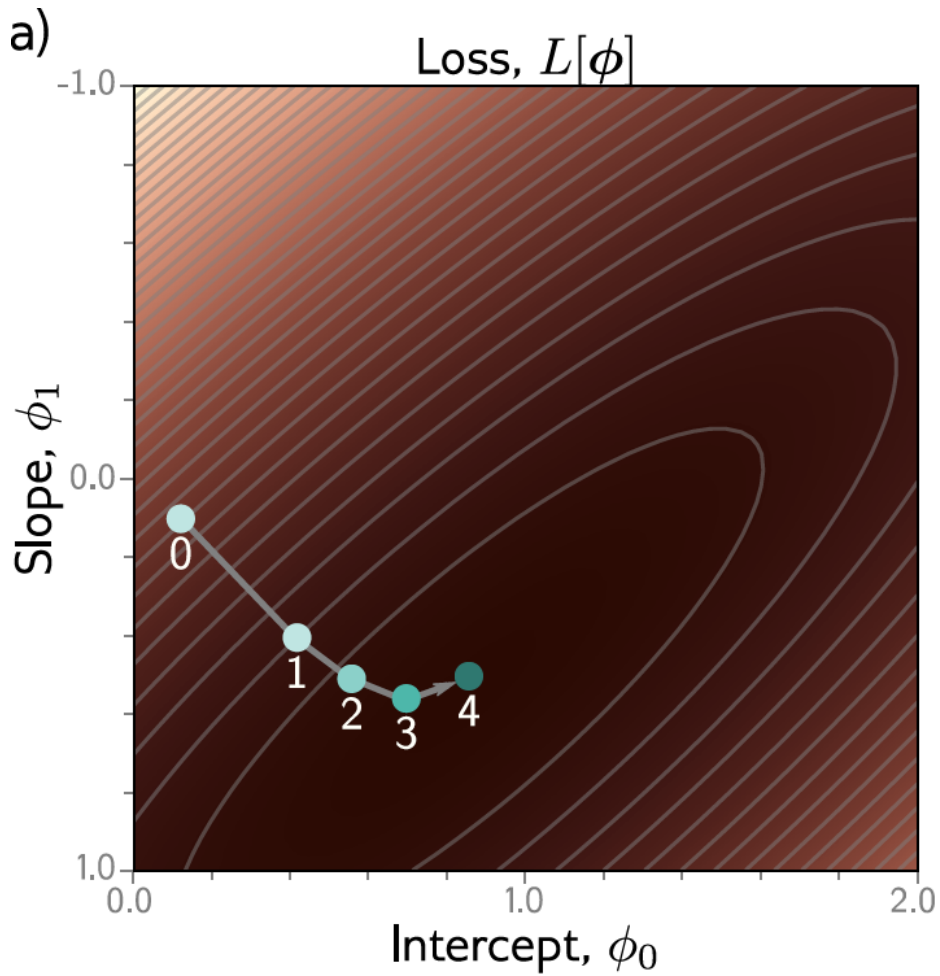- Fully connected network

# Recap: 1D Linear regression loss function



Loss function:

$$L[\boldsymbol{\phi}] = \sum_{i=1}^{I} (\mathrm{f}[x_i, \boldsymbol{\phi}] - y_i)^2$$

$$= \sum_{i=1}^{I} (\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"

# Recap: Gradient Descent

a)

Loss, $L[\phi]$

b)

# Shallow neural networks

- 1D regression model is obviously limited
  - Want to be able to describe input/output that are not lines
  - Want multiple inputs
  - Want multiple outputs

- Shallow neural networks
  - Flexible enough to describe arbitrarily complex input/output mappings
  - Can have as many inputs as we want
  - Can have as many outputs as we want

# Shallow Neural Networks

- Example network, 1 input, 1 output
- Universal approximation theorem
- More than one output
- More than one input
- General case
- Number of regions
- Terminology

# 1D Linear Regression

$$y = f[x, \phi]$$

$$= \phi_0 + \phi_1 x$$

# Example shallow network

$$y = f[x, \phi]$$

$$= \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$$

# Example shallow network

generic interface ✓

many $\phi_i$'s

$$y = f[x, \phi]$$

$$= \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$$

# Example shallow network

Activation function

$$y = \mathrm{f}[x, \phi]$$

$$= \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$$
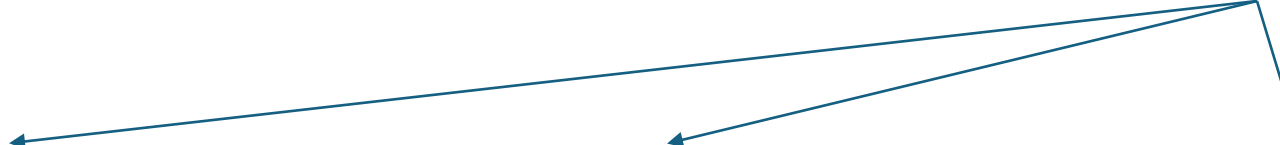
linear func     linear     linear

intermediate outputs

output is linear combination of

whole function linear except activation function

# Example shallow network

Activation function

$$y = \mathrm{f}[x, \phi]$$

$$= \phi_0 + \phi_1 \mathrm{a}[\theta_{10} + \theta_{11}x] + \phi_2 \mathrm{a}[\theta_{20} + \theta_{21}x] + \phi_3 \mathrm{a}[\theta_{30} + \theta_{31}x]$$

$$\mathrm{a}[z] = \mathrm{ReLU}[z] = \begin{cases} 0 & z < 0 \\ z & z \geq 0 \end{cases}.$$
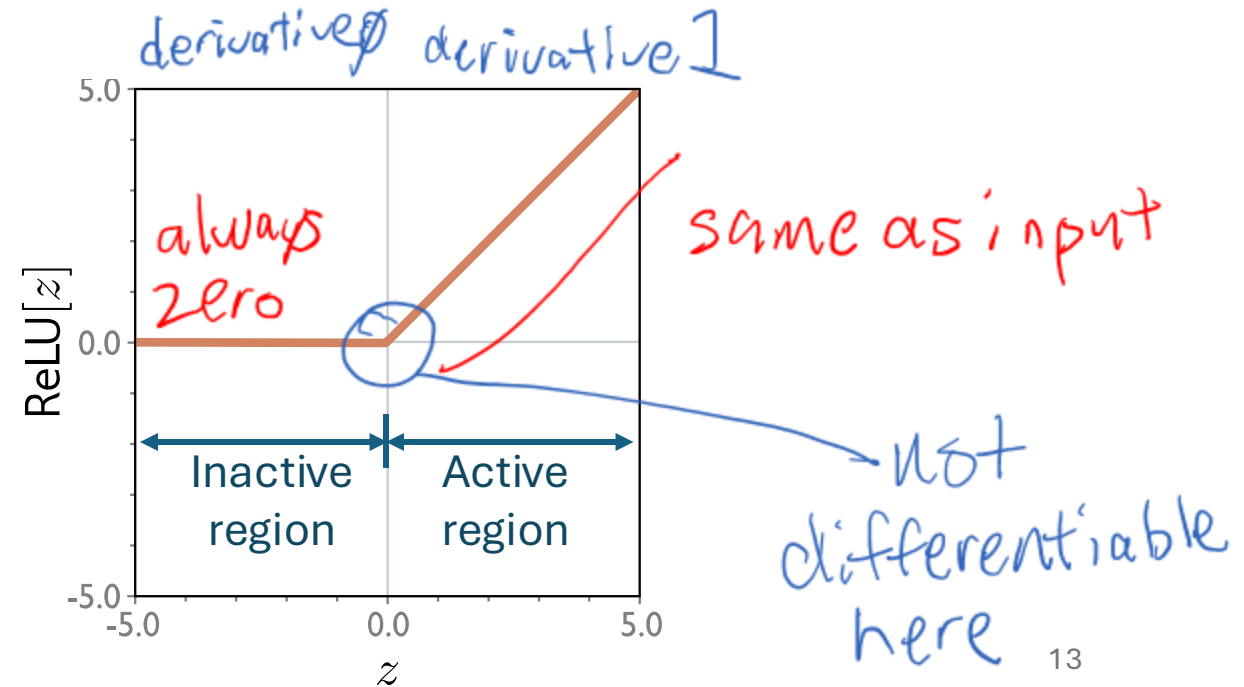
Rectified Linear Unit

(one type of activation function)

# Example shallow network

$$y = \mathrm{f}[x, \phi]$$

$$= \phi_0 + \phi_1 \mathrm{a}[\theta_{10} + \theta_{11}x] + \phi_2 \mathrm{a}[\theta_{20} + \theta_{21}x] + \phi_3 \mathrm{a}[\theta_{30} + \theta_{31}x]$$

Activation function

$$\mathrm{a}[z] = \mathrm{ReLU}[z] = \begin{cases} 0 & z < 0 \\ z & z \geq 0 \end{cases}.$$

Rectified Linear Unit
(particular kind of activation function)



derivative 0   derivative 1

always zero

same as input

Inactive region    Active region

not differentiable here

13

# Example shallow network

*(handwritten annotations: +1, x3, 3 params)*

$$y = \mathrm{f}[x, \phi]$$

$$= \phi_0 + \phi_1 \mathrm{a}[\theta_{10} + \theta_{11}x] + \phi_2 \mathrm{a}[\theta_{20} + \theta_{21}x] + \phi_3 \mathrm{a}[\theta_{30} + \theta_{31}x]$$

This model has 10 parameters:

$$\phi = \{\phi_0, \phi_1, \phi_2, \phi_3, \theta_{10}, \theta_{11}, \theta_{20}, \theta_{21}, \theta_{30}, \theta_{31}\}$$

- Represents a family of functions
- Parameters determine a particular function
- Given the parameters, we can perform inference (evaluate the equation) $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{I}$
- Given training dataset $L[\phi]$
- Define loss function (least squares)
- Change parameters to minimize loss function
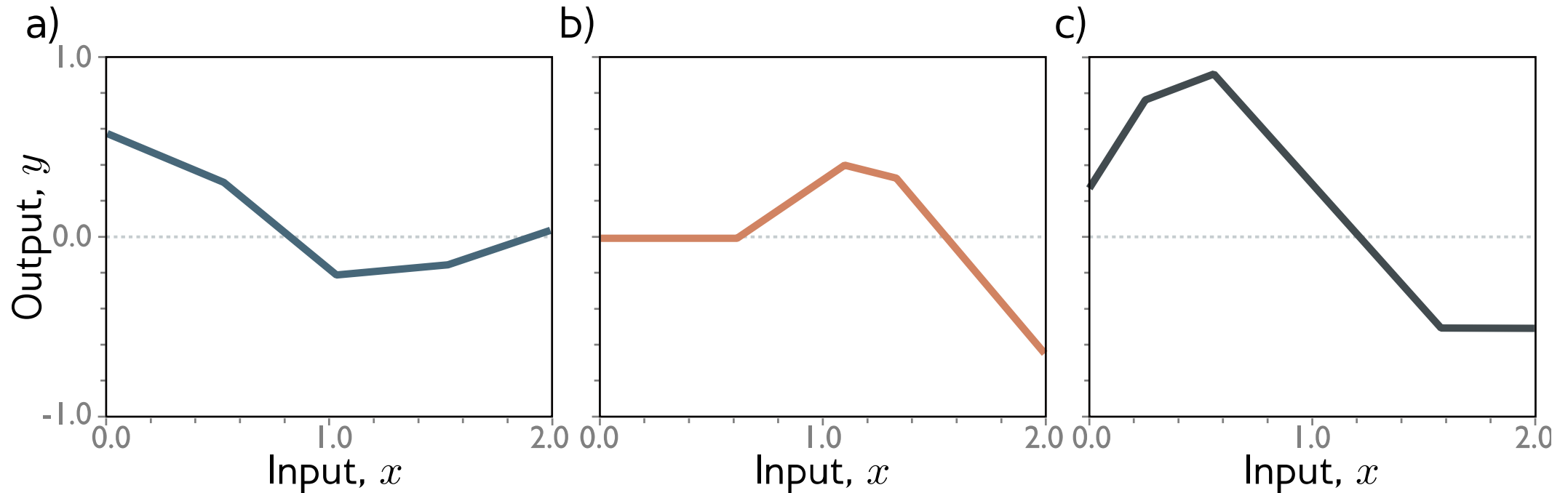
# Example shallow network

$$y = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x].$$

# Example shallow network

3 joints from
3 activation functions

0 in a[] input will be joint

$$y = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x].$$

a)

b)

c)



Piecewise linear functions with three joints

16

# Hidden units

$$y = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x].$$
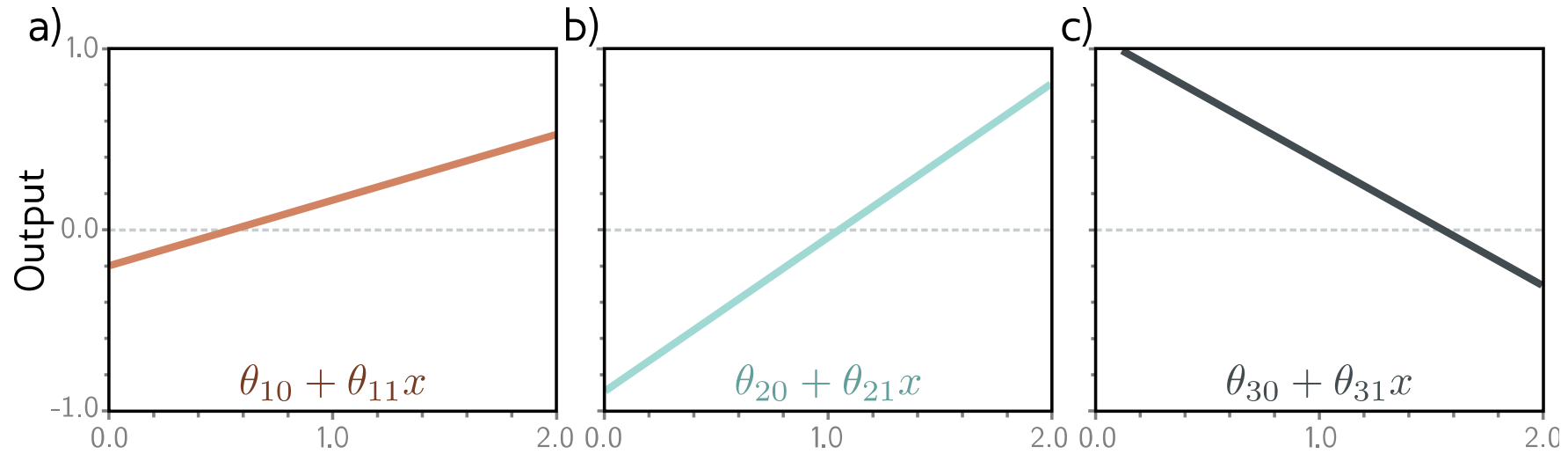
Break down into two parts:

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

where:

Hidden units
$$h_1 = a[\theta_{10} + \theta_{11}x]$$
$$h_2 = a[\theta_{20} + \theta_{21}x]$$
$$h_3 = a[\theta_{30} + \theta_{31}x]$$

# 1. compute three linear functions

**Linear Functions**

a)

$\theta_{10} + \theta_{11}x$

b)

$\theta_{20} + \theta_{21}x$

c)

$\theta_{30} + \theta_{31}x$

Output

1.0

0.0

-1.0

0.0    1.0    2.0

0.0    1.0    2.0
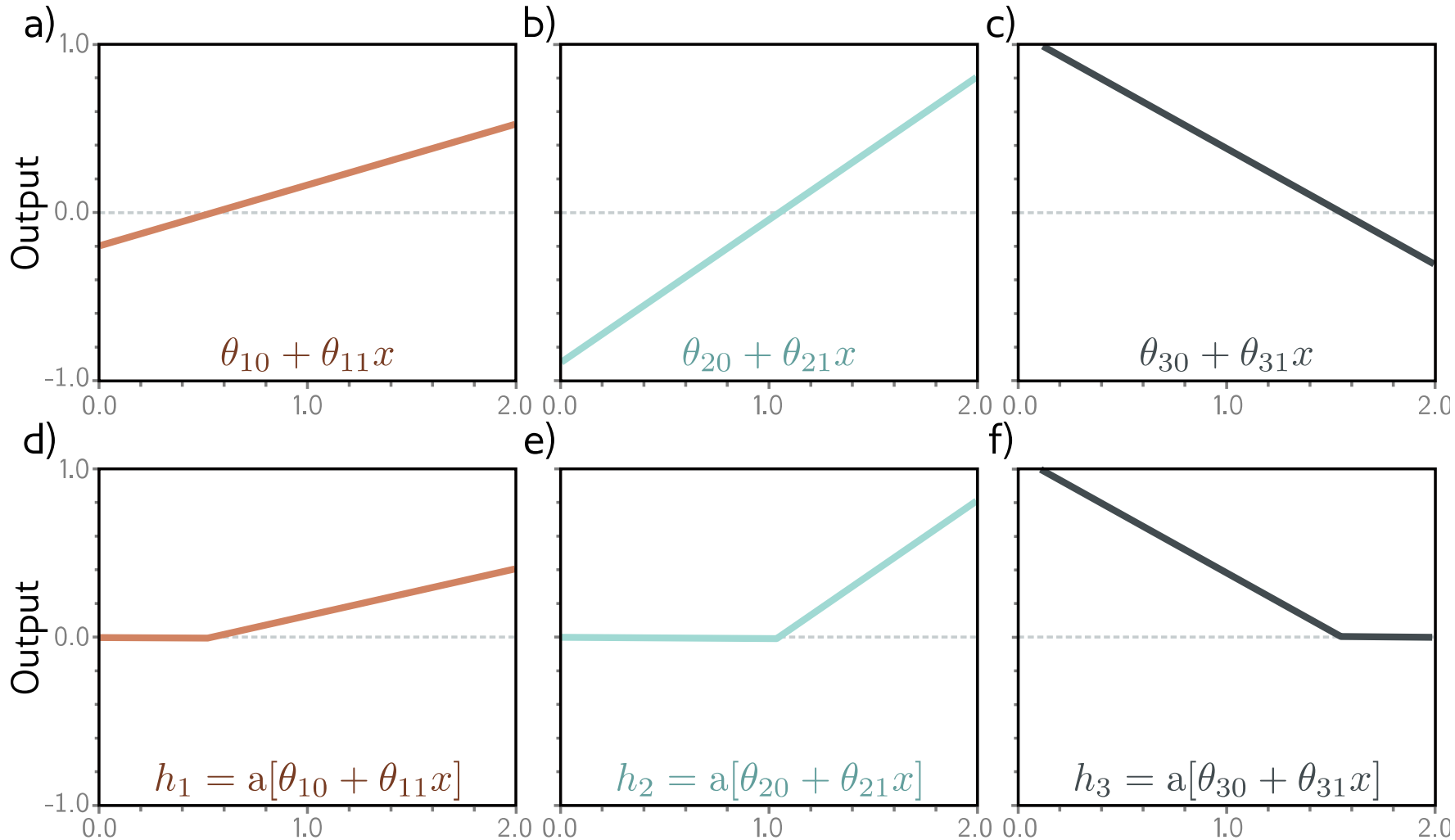
0.0    1.0    2.0

## 2. Pass through ReLU functions (creates hidden units)

$$h_1 = \text{a}[\theta_{10} + \theta_{11}x]$$
$$h_2 = \text{a}[\theta_{20} + \theta_{21}x]$$
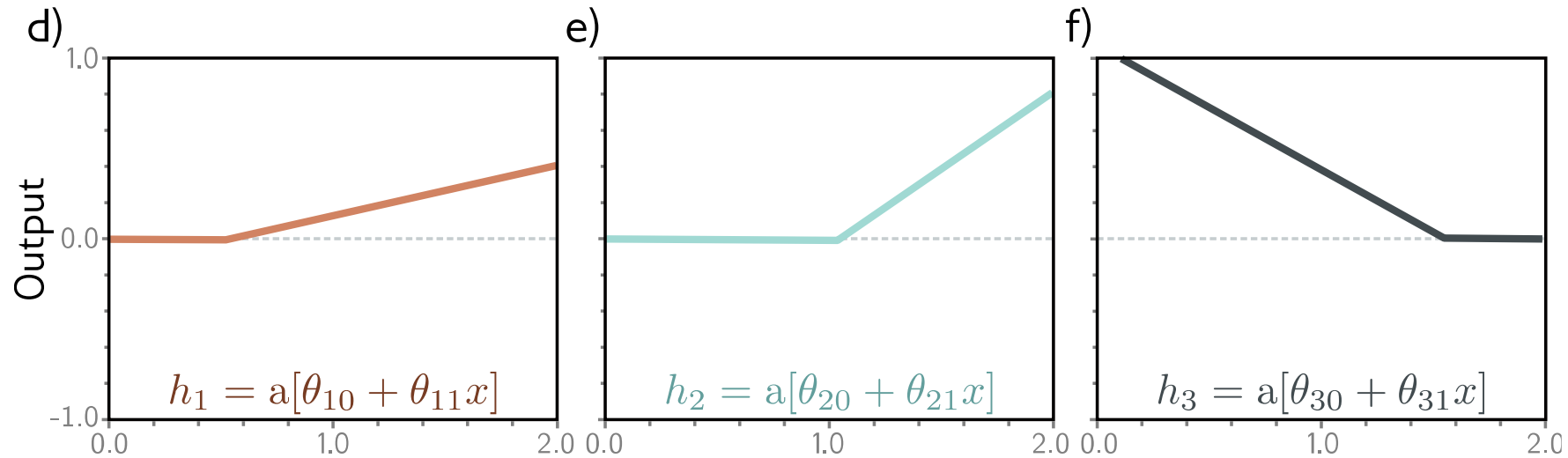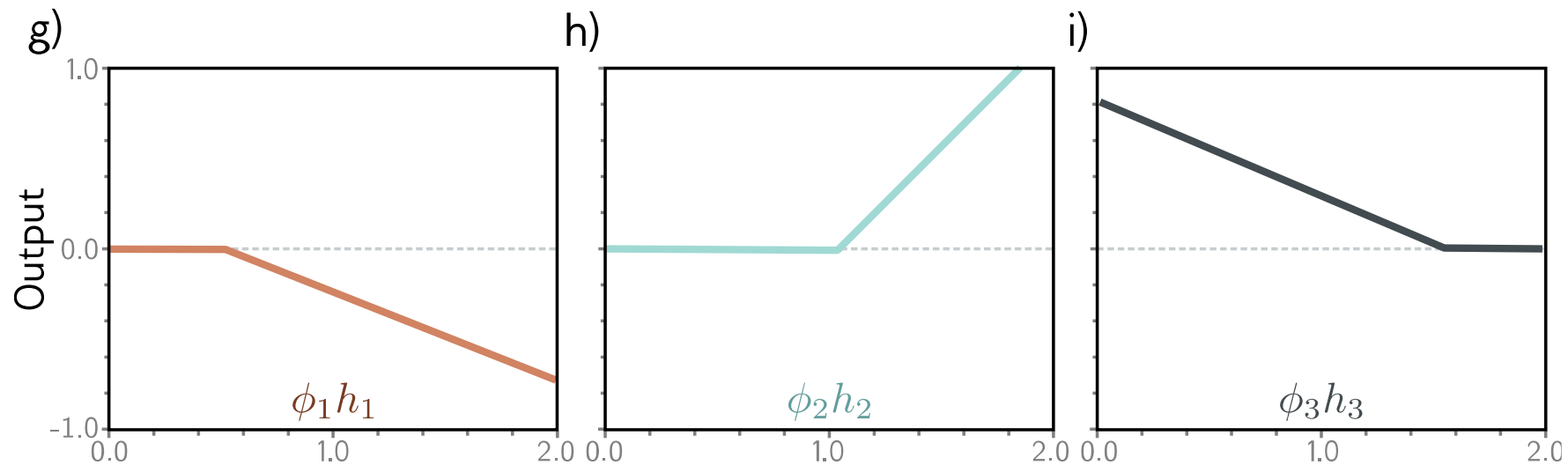$$h_3 = \text{a}[\theta_{30} + \theta_{31}x],$$

**Linear Functions**

a)

$\theta_{10} + \theta_{11}x$

b)

$\theta_{20} + \theta_{21}x$

c)

$\theta_{30} + \theta_{31}x$

**After Activation**

d)

$h_1 = \text{a}[\theta_{10} + \theta_{11}x]$

e)

$h_2 = \text{a}[\theta_{20} + \theta_{21}x]$

f)

$h_3 = \text{a}[\theta_{30} + \theta_{31}x]$

# 2. Weight the hidden units

**After Activation**

d)

$h_1 = \mathrm{a}[\theta_{10} + \theta_{11}x]$

e)

$h_2 = \mathrm{a}[\theta_{20} + \theta_{21}x]$

f)

$h_3 = \mathrm{a}[\theta_{30} + \theta_{31}x]$

**Weight the Hidden units**

g)

$\phi_1 h_1$

h)

$\phi_2 h_2$

i)

$\phi_3 h_3$

# 4. Sum the weighted hidden units to create output

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

**Weight the hidden units**

g)

$\phi_1 h_1$

Output

Input, $x$

h)

$\phi_2 h_2$

i)

$\phi_3 h_3$

Input, $x$

**Sum the weighted hidden units**

j)

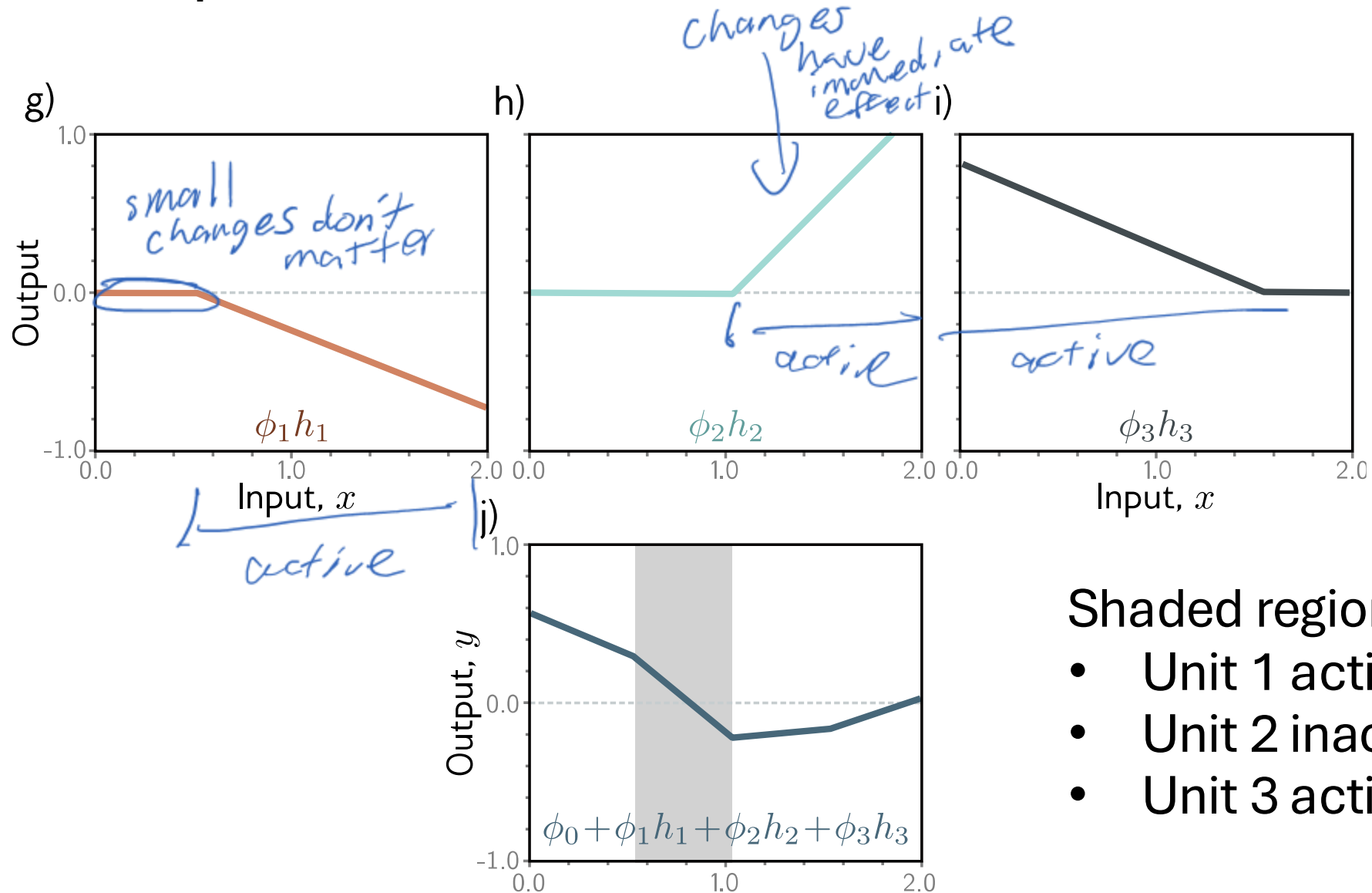$\phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$

Output, $y$

Input, $x$

# Example: 3 different shallow networks

$$y = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x].$$
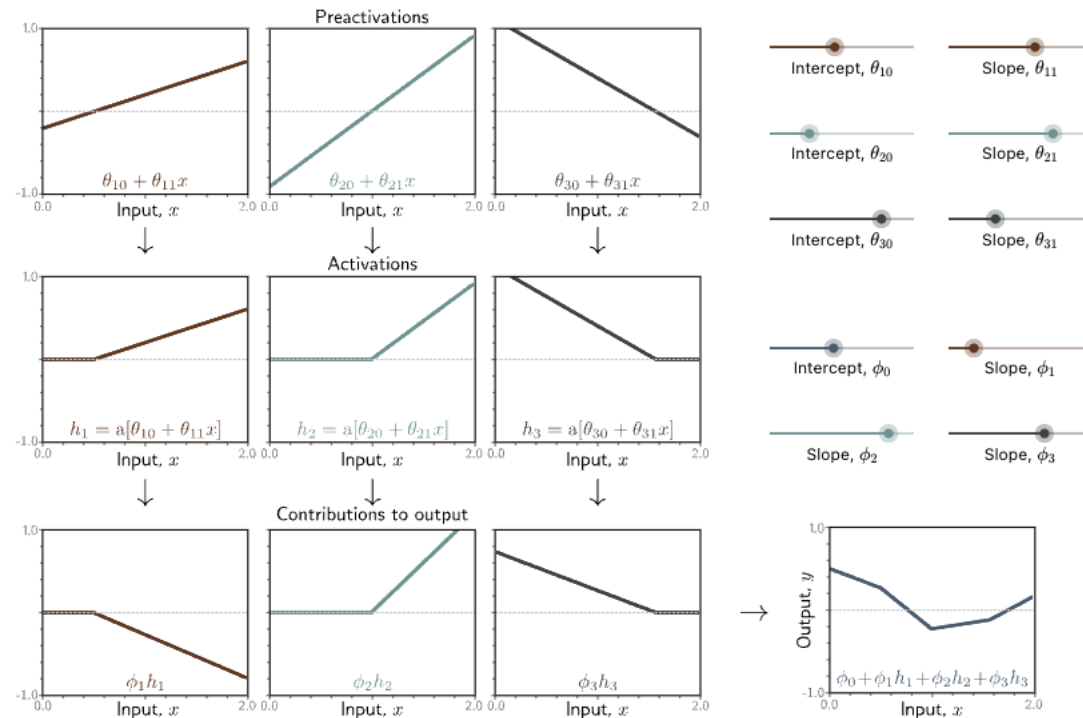


Example shallow network = piecewise linear functions
1 "joint" per ReLU function

# Activation pattern = which hidden units are activated?

g)



$\phi_1 h_1$

small changes don't matter

active

h)



$\phi_2 h_2$

changes have immediate effect

active

i)



$\phi_3 h_3$

active

j)



$\phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$

Shaded region:
- Unit 1 active
- Unit 2 inactive
- Unit 3 active

# Interactive Figure 3.3a: 1D Shallow Network (ReLU)



**Figure 3.3** Computation for function in figure 3.2a. (Top row) The input $x$ is passed through three linear functions, each with a different y-intercept $\theta_{\bullet 0}$ and slope $\theta_{\bullet 1}$. (Center row) Each line is passed through the ReLU activation function. (Bottom row) The three resulting functions are then weighted (scaled) by $\phi_1, \phi_2$, and $\phi_3$, respectively. (Bottom right) Finally, the weighted functions are summed, and an offset $\phi_0$ that controls the height is added.

Move the sliders to modify the parameters of the shallow network.

https://udlbook.github.io/udlfigures/

# Depicting neural networks

hidden units

$$h_1 = \mathrm{a}[\theta_{10} + \theta_{11}x]$$

$$h_2 = \mathrm{a}[\theta_{20} + \theta_{21}x]$$

$$h_3 = \mathrm{a}[\theta_{30} + \theta_{31}x]$$

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$



Each parameter multiplies its source and adds to its target
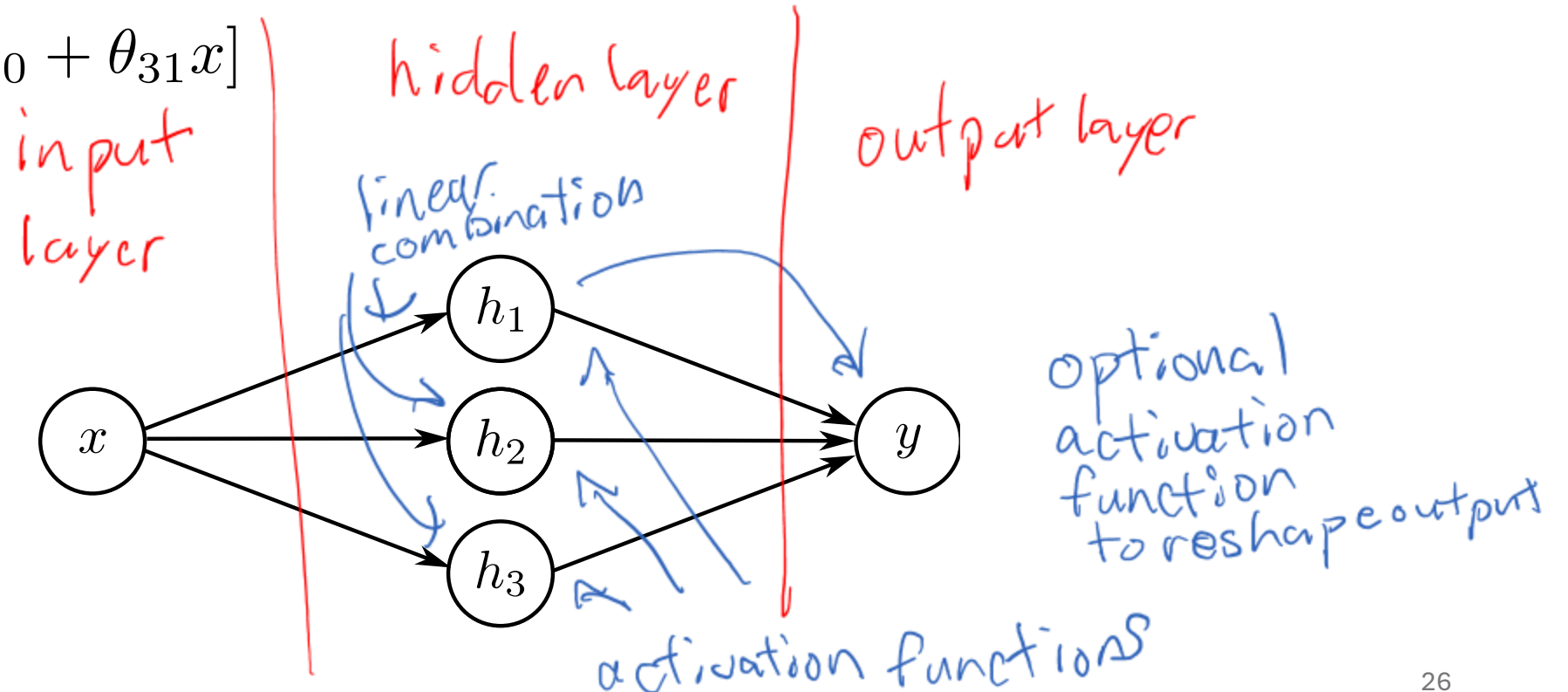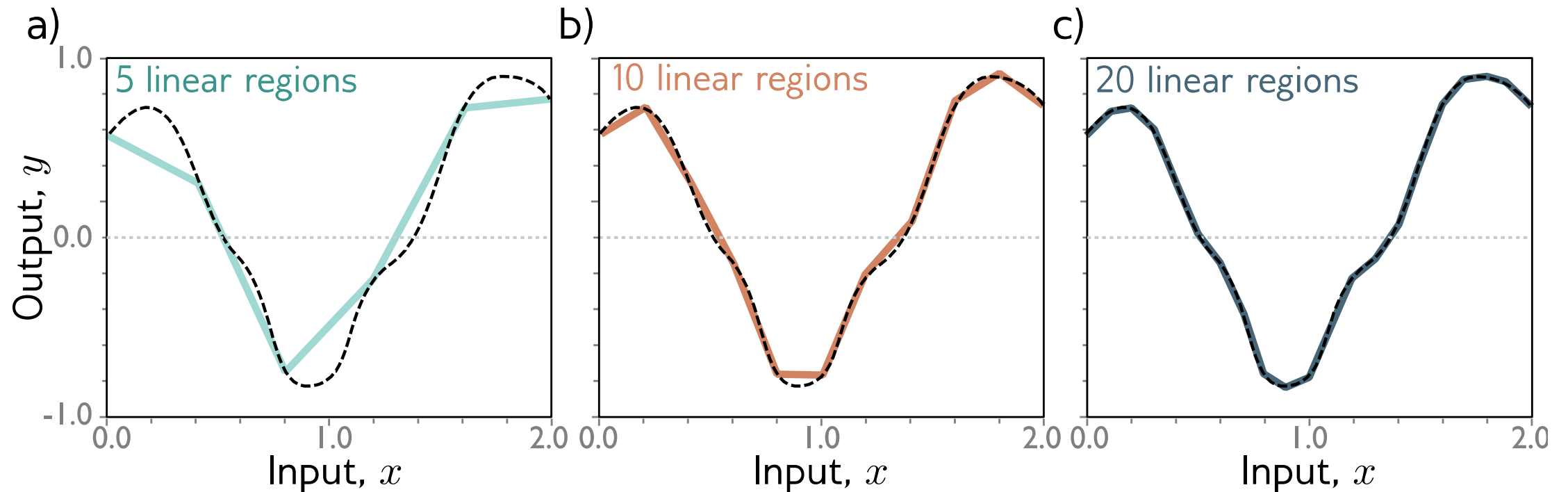
# Depicting neural networks
## *Usually don't show the bias terms*

$$h_1 = \mathrm{a}[\theta_{10} + \theta_{11}x]$$

$$h_2 = \mathrm{a}[\theta_{20} + \theta_{21}x] \qquad y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

$$h_3 = \mathrm{a}[\theta_{30} + \theta_{31}x]$$

hidden layer

input layer

linear combination

output layer

optional activation function to reshape output

$x$ → $h_1$, $h_2$, $h_3$ → $y$

activation functions

# Any questions?

# Shallow neural networks

- Example network, 1 input, 1 output
- Universal approximation theorem
- More than one output
- More than one input
- General case
- Number of regions
- Terminology

# With 3 hidden units:

$$h_1 = \mathrm{a}[\theta_{10} + \theta_{11}x]$$
$$h_2 = \mathrm{a}[\theta_{20} + \theta_{21}x]$$
$$h_3 = \mathrm{a}[\theta_{30} + \theta_{31}x]$$

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

# With D hidden units:

$$h_d = \mathrm{a}[\theta_{d0} + \theta_{d1}x]$$

$$y = \phi_0 + \sum_{d=1}^{D} \phi_d h_d$$

# With enough hidden units...

... we can describe any 1D function to arbitrary accuracy



a) 5 linear regions

b) 10 linear regions

c) 20 linear regions

# Universal approximation theorem*s*

"a formal proof that, with enough hidden units, a shallow neural network can describe any continuous function in $R^D$ to arbitrary precision"

#of hidden layers → $\emptyset$ does not work
1+ does

activation functions → linear does not work
Polynomials do not work

1+ hidden and non-polynomial activation
works

# Any questions?

# Shallow neural networks

- Example network, 1 input, 1 output
- Universal approximation theorem
- More than one output
- More than one input
- General case
- Number of regions
- Terminology

# Two outputs

- 1 input, 4 hidden units, 2 outputs

$$h_1 = a[\theta_{10} + \theta_{11}x]$$
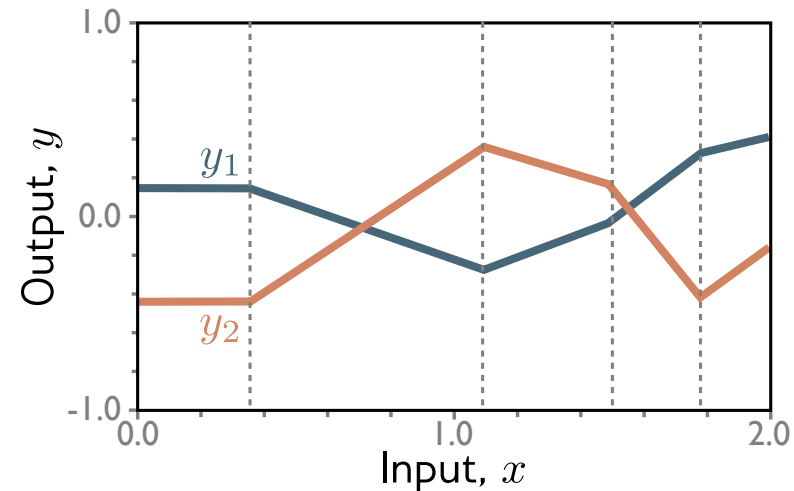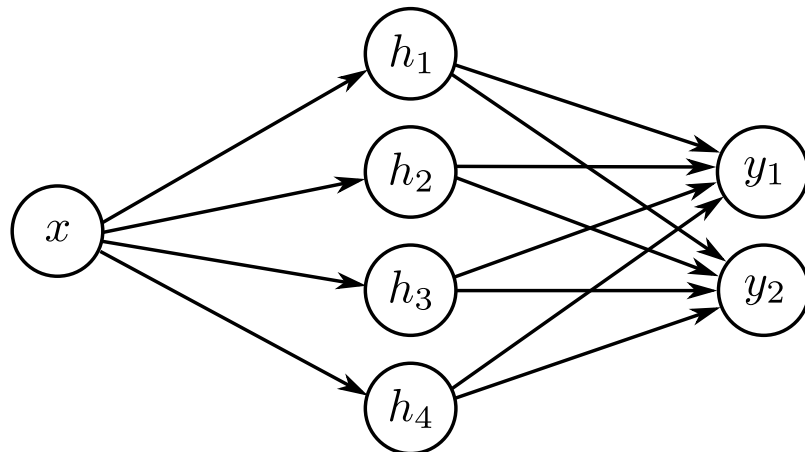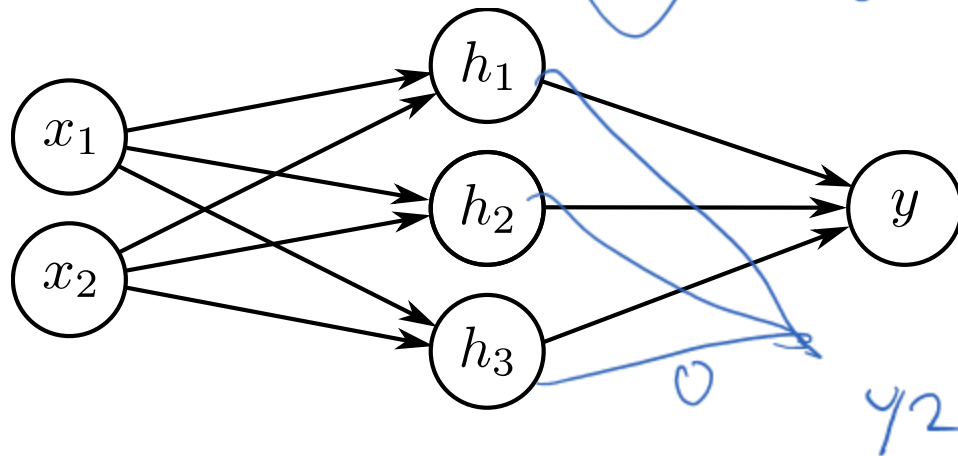$$h_2 = a[\theta_{20} + \theta_{21}x]$$
$$h_3 = a[\theta_{30} + \theta_{31}x]$$
$$h_4 = a[\theta_{40} + \theta_{41}x]$$

*same as before*

$$y_1 = \phi_{10} + \phi_{11}h_1 + \phi_{12}h_2 + \phi_{13}h_3 + \phi_{14}h_4$$
$$y_2 = \phi_{20} + \phi_{21}h_1 + \phi_{22}h_2 + \phi_{23}h_3 + \phi_{24}h_4$$

*two output nodes + formulas*

# Two outputs

- 1 input, 4 hidden units, 2 outputs

$$h_1 = \mathrm{a}[\theta_{10} + \theta_{11}x]$$
$$h_2 = \mathrm{a}[\theta_{20} + \theta_{21}x]$$
$$h_3 = \mathrm{a}[\theta_{30} + \theta_{31}x]$$
$$h_4 = \mathrm{a}[\theta_{40} + \theta_{41}x]$$

$$y_1 = \phi_{10} + \phi_{11}h_1 + \phi_{12}h_2 + \phi_{13}h_3 + \phi_{14}h_4$$
$$y_2 = \phi_{20} + \phi_{21}h_1 + \phi_{22}h_2 + \phi_{23}h_3 + \phi_{24}h_4$$



one more output node.
same structure, different parameters.

# Two outputs

- 1 input, 4 hidden units, 2 outputs

$$h_1 = a[\theta_{10} + \theta_{11}x]$$
$$h_2 = a[\theta_{20} + \theta_{21}x]$$
$$h_3 = a[\theta_{30} + \theta_{31}x]$$
$$h_4 = a[\theta_{40} + \theta_{41}x]$$

$$y_1 = \phi_{10} + \phi_{11}h_1 + \phi_{12}h_2 + \phi_{13}h_3 + \phi_{14}h_4$$
$$y_2 = \phi_{20} + \phi_{21}h_1 + \phi_{22}h_2 + \phi_{23}h_3 + \phi_{24}h_4$$

joints @ same inputs

# Any questions?

# Shallow neural networks

- Example network, 1 input, 1 output
- Universal approximation theorem
- More than one output
- More than one input
- General case
- Number of regions
- Terminology

# Two inputs

- 2 inputs, 3 hidden units, 1 output

*1 input*

$h_1 = \mathrm{a}[\theta_{10} + \theta_{11}x_1 + \theta_{12}x_2]$

$h_2 = \mathrm{a}[\theta_{20} + \theta_{21}x_1 + \theta_{22}x_2]$

$h_3 = \mathrm{a}[\theta_{30} + \theta_{31}x_1 + \theta_{32}x_2]$

*new*

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

**Linear Functions**

Plots titled (left to right): $\theta_{10} + \theta_{11}x_1 + \theta_{12}x_2$, $\theta_{20} + \theta_{21}x_1 + \theta_{22}x_2$, $\theta_{30} + \theta_{31}x_1 + \theta_{32}x_2$. Each plot has axes Input, $x_1$ (horizontal) and Input, $x_2$ (vertical), ranging from -1.0 to 1.0.

See Interactive Figure 3.8a https://udlbook.github.io/udlfigures/

**Linear Functions**

$$\theta_{10} + \theta_{11}x_1 + \theta_{12}x_2$$

$$\theta_{20} + \theta_{21}x_1 + \theta_{22}x_2$$

$$\theta_{30} + \theta_{31}x_1 + \theta_{32}x_2$$

*negative*

*positive*

**After Activation**

$$h_1 = a[\theta_{10} + \theta_{11}x_1 + \theta_{12}x_2]$$

$$h_2 = a[\theta_{20} + \theta_{21}x_1 + \theta_{22}x_2]$$

$$h_3 = a[\theta_{30} + \theta_{31}x_1 + \theta_{32}x_2]$$

*zero, previously neg*

Input, $x_1$   Input, $x_2$

41

After Activation

$h_1 = \mathrm{a}[\theta_{10} + \theta_{11}x_1 + \theta_{12}x_2]$

$h_2 = \mathrm{a}[\theta_{20} + \theta_{21}x_1 + \theta_{22}x_2]$

$h_3 = \mathrm{a}[\theta_{30} + \theta_{31}x_1 + \theta_{32}x_2]$

Weight the Hidden units

$\phi_1 h_1$

$\phi_2 h_2$

$\phi_3 h_3$

42

$\phi_1 h_1$     $\phi_2 h_2$     $\phi_3 h_3$

**Weight the hidden units**

**Sum the weighted hidden units**

$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$

43

$\phi_1 h_1$  $\phi_2 h_2$  $\phi_3 h_3$

$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$

Interactive Figure 3.8b

https://udlbook.github.io/udlfigures/

# Convex polygonal regions

A region of $\mathbb{R}^D$ is convex if we can draw a straight line between any two points on the boundary of the region without intersecting the boundary in another place.

# Fitting a dataset where:
## each sample has 2 inputs and 1 output



$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

# Question:

- For the 2D case, what if there were two outputs?

- If this is one of the outputs, what would the other one look like?



$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

# Any questions?

# Shallow neural networks

- Example network, 1 input, 1 ouput
- Universal approximation theorem
- More than one output
- More than one input
- General case
- Number of regions
- Terminology

# Arbitrary inputs, hidden units, outputs

- $D_i$ inputs, $D$ hidden units, and $D_o$ Outputs

$$h_d = \mathrm{a}\left[\theta_{d0} + \sum_{i=1}^{D_i} \theta_{di} x_i\right] \qquad y_j = \phi_{j0} + \sum_{d=1}^{D} \phi_{jd} h_d$$
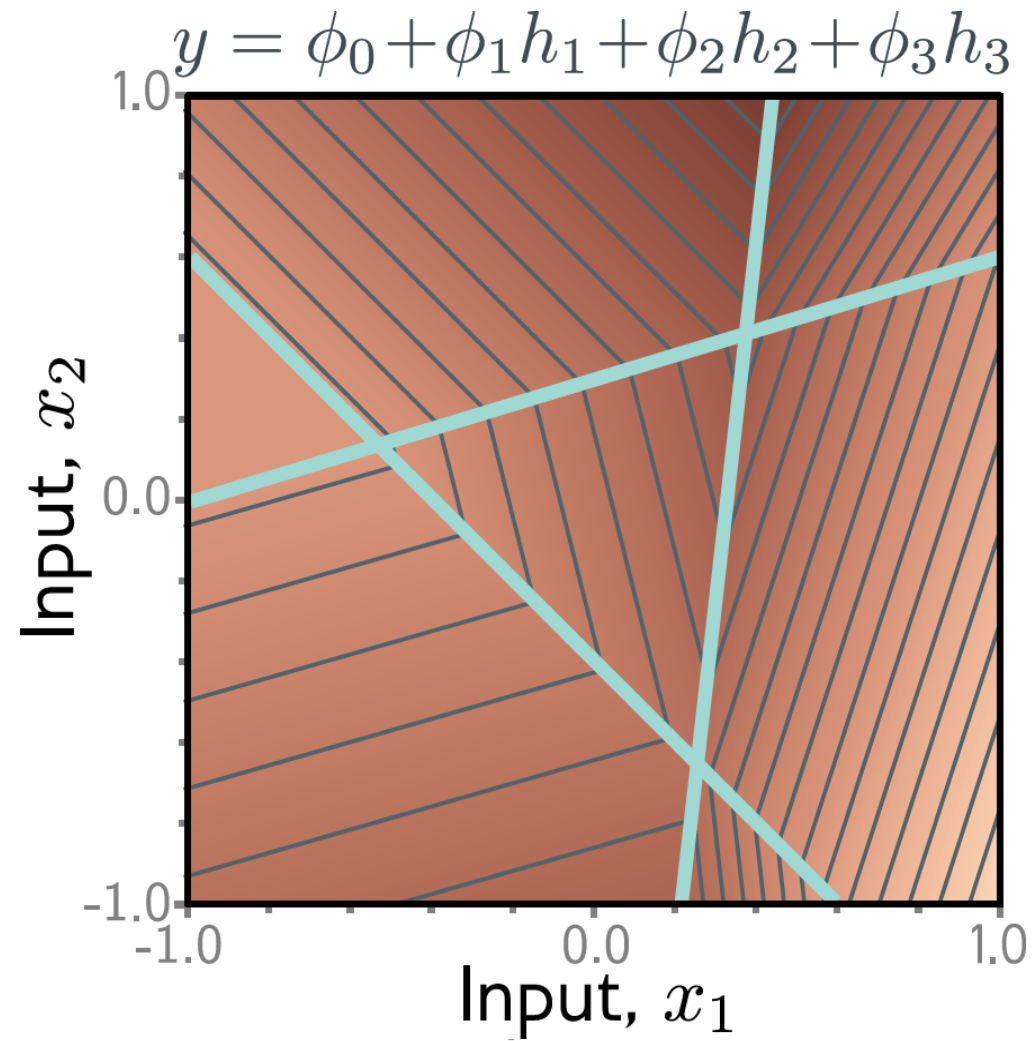
- e.g., Three inputs, three hidden units, two outputs

# Question:

- How many parameters does this model have?

# How many hidden units?

Input

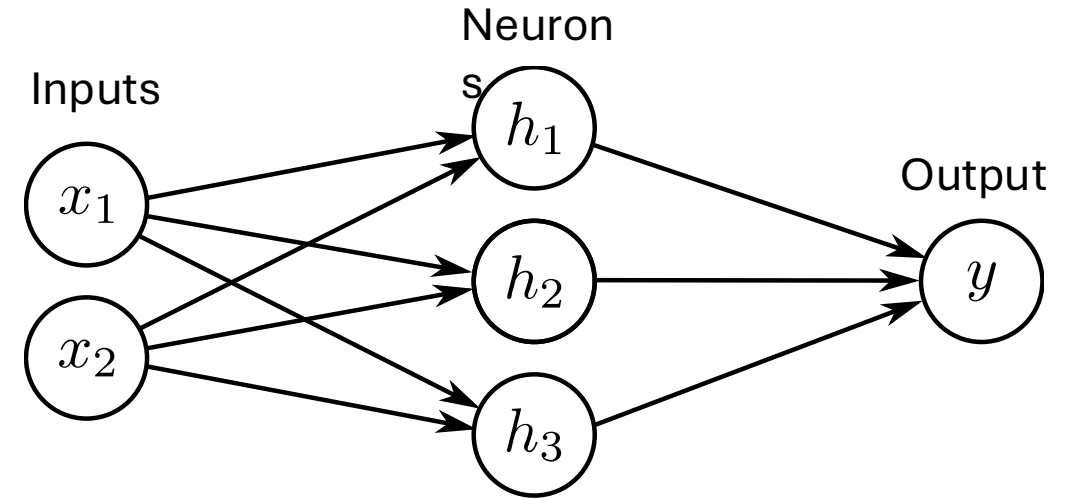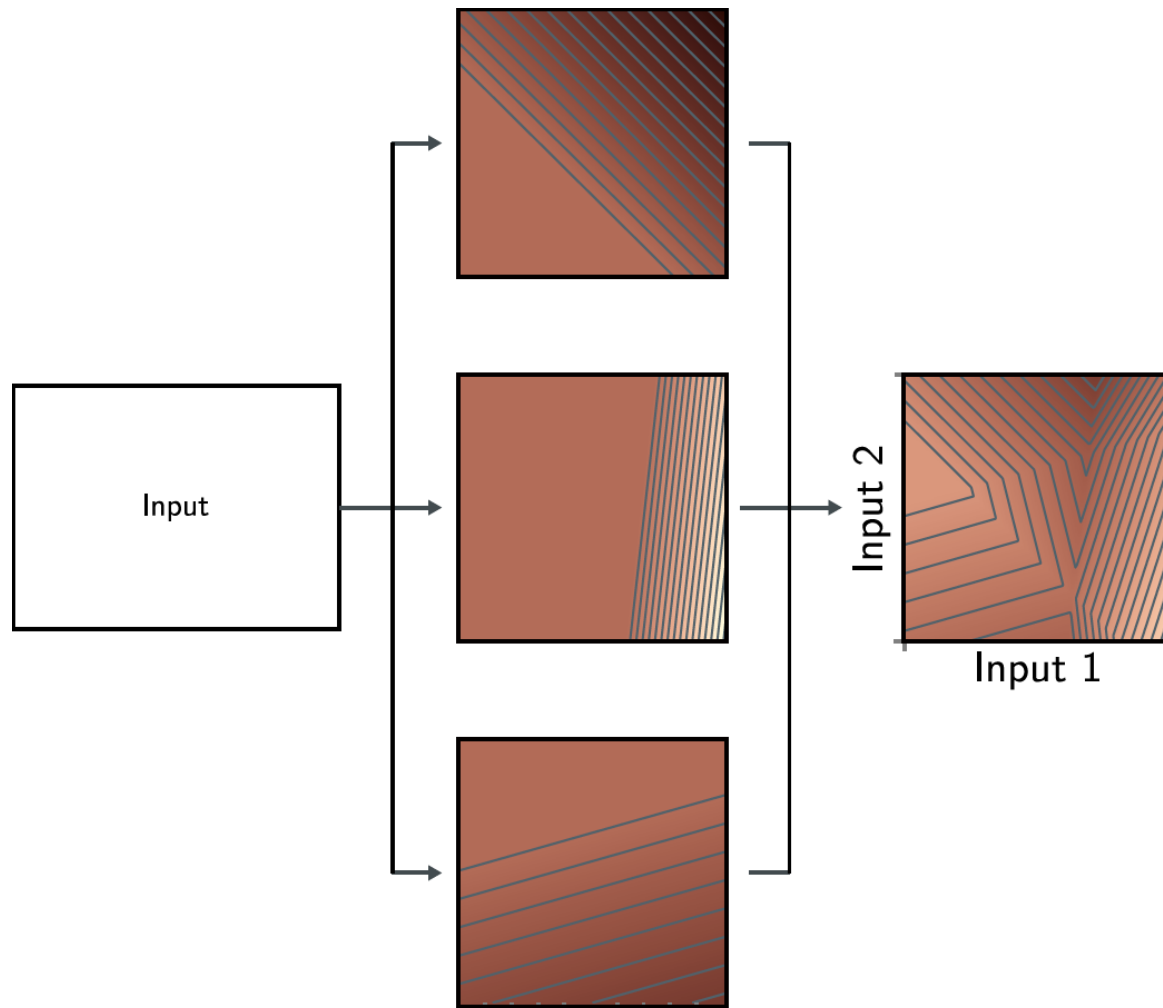# Output with boundaries and in 3D

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

# How would you draw and write this neural network?



Input

Input 2

Input 1

# How would you draw and write this neural network?

Inputs

Neurons

$h_1$

$h_2$

$h_3$

Output

$x_1$

$x_2$

$y$

"neural network"

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

$$h_1 = a[\theta_{10} + \theta_{11} x_1 + \theta_{12} x_2]$$
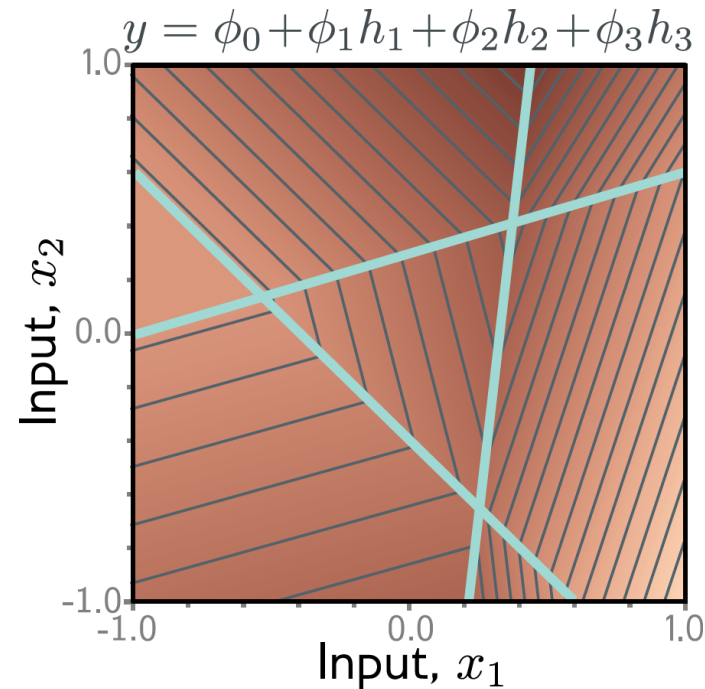$$h_2 = a[\theta_{20} + \theta_{21} x_1 + \theta_{22} x_2]$$
$$h_3 = a[\theta_{30} + \theta_{31} x_1 + \theta_{32} x_2]$$

Input

Input 2

Input 1

# Shallow neural networks

- Example network, 1 input, 1 output
- Universal approximation theorem
- More than one output
- More than one input
- General case
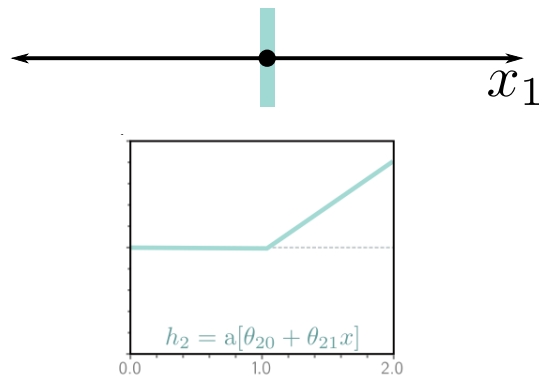- Number of regions
- Terminology

# Number of output regions

- With ReLU activations, each output consists of multi-dimensional piecewise linear hyperplanes

- With two inputs, and three hidden units, we saw there were seven polygons for each output:

# Example with $D = D_i \rightarrow 2^{D_i}$ regions

a)

b)

c)



$h_2 = \text{a}[\theta_{20} + \theta_{21}x]$
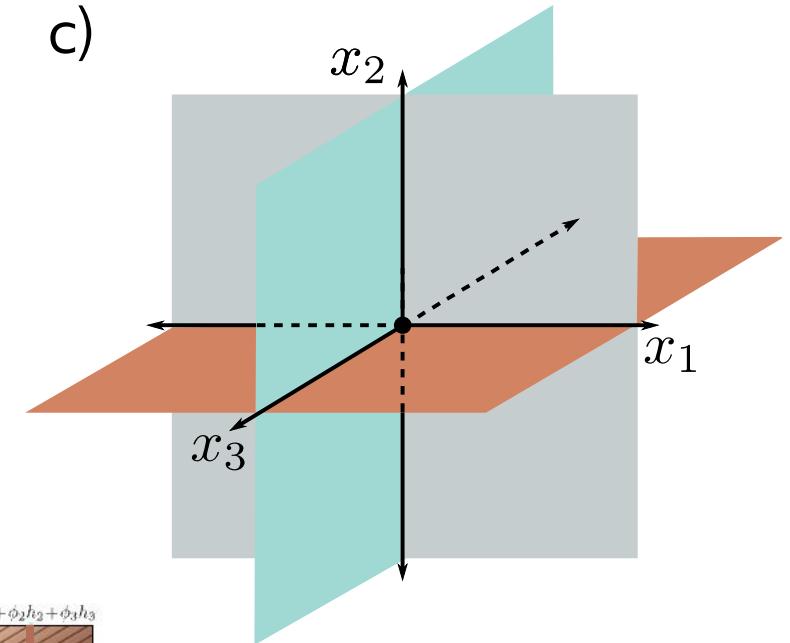
- 1 input (1-dimension)
- 1 hidden unit
- creates two regions (one joint)

- 2 input (2-dimensions) with
- 2 hidden units
- creates four regions (two lines)

- 3 inputs (3-dimensions) with
- 3 hidden units
- creates eight regions (three planes)

58

# Number of regions:

$D_i$ : # of inputs
$D$ : # of hidden units
$D_o$ : # of outputs

*# of hidden units* (handwritten)

- Number of regions created by $D > D_i$ hyper-planes in $D_i$ *# inputs* (handwritten) dimensions was proved by Zaslavsky (1975) to be:

$$\sum_{j=0}^{D_i} \binom{D}{j} = \frac{D!}{D_i!(D-D_i)!}$$

←———— Binomial coefficients!

- How big is this?  It's greater than $2^{D_i}$ but less than $2^{D}$.

# Number of output regions

$D_i$ : # of inputs
$D$ : # of hidden units
$D_o$ : # of outputs

- In general, each output consists of D dimensional convex polytopes
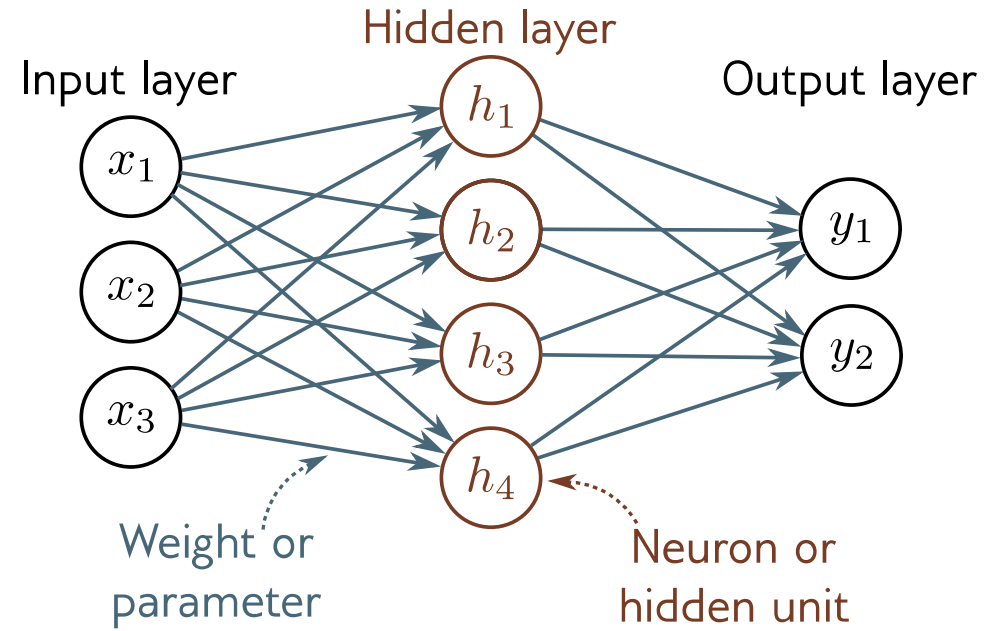
- How many?

a)



b)

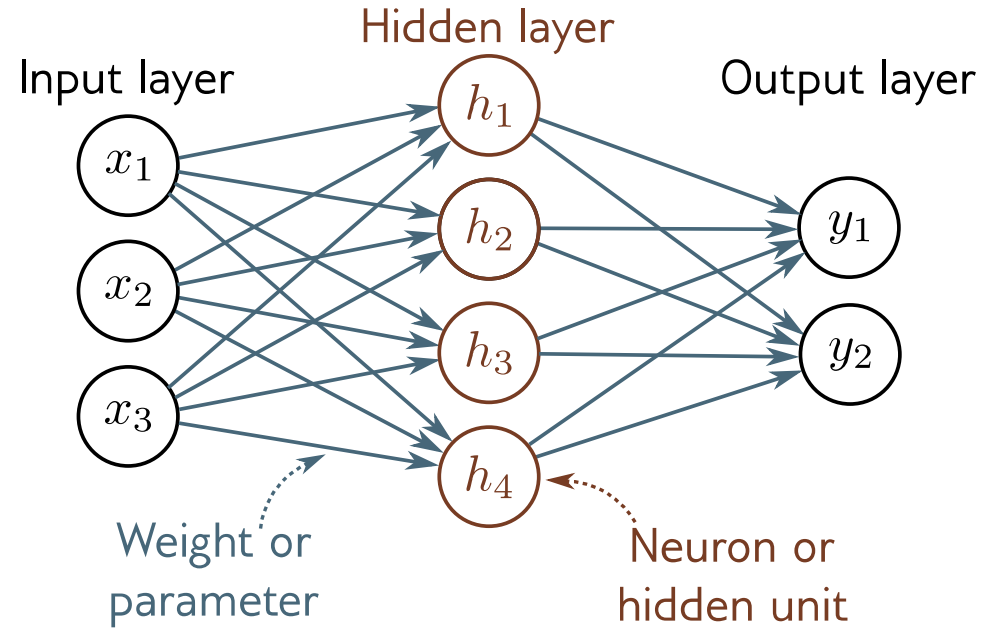Highlighted point = 500 hidden units or 51,001 parameters

60

# Any questions?

# Shallow neural networks

- Example network, 1 input, 1 output
- Universal approximation theorem
- More than one output
- More than one input
- General case
- Number of regions
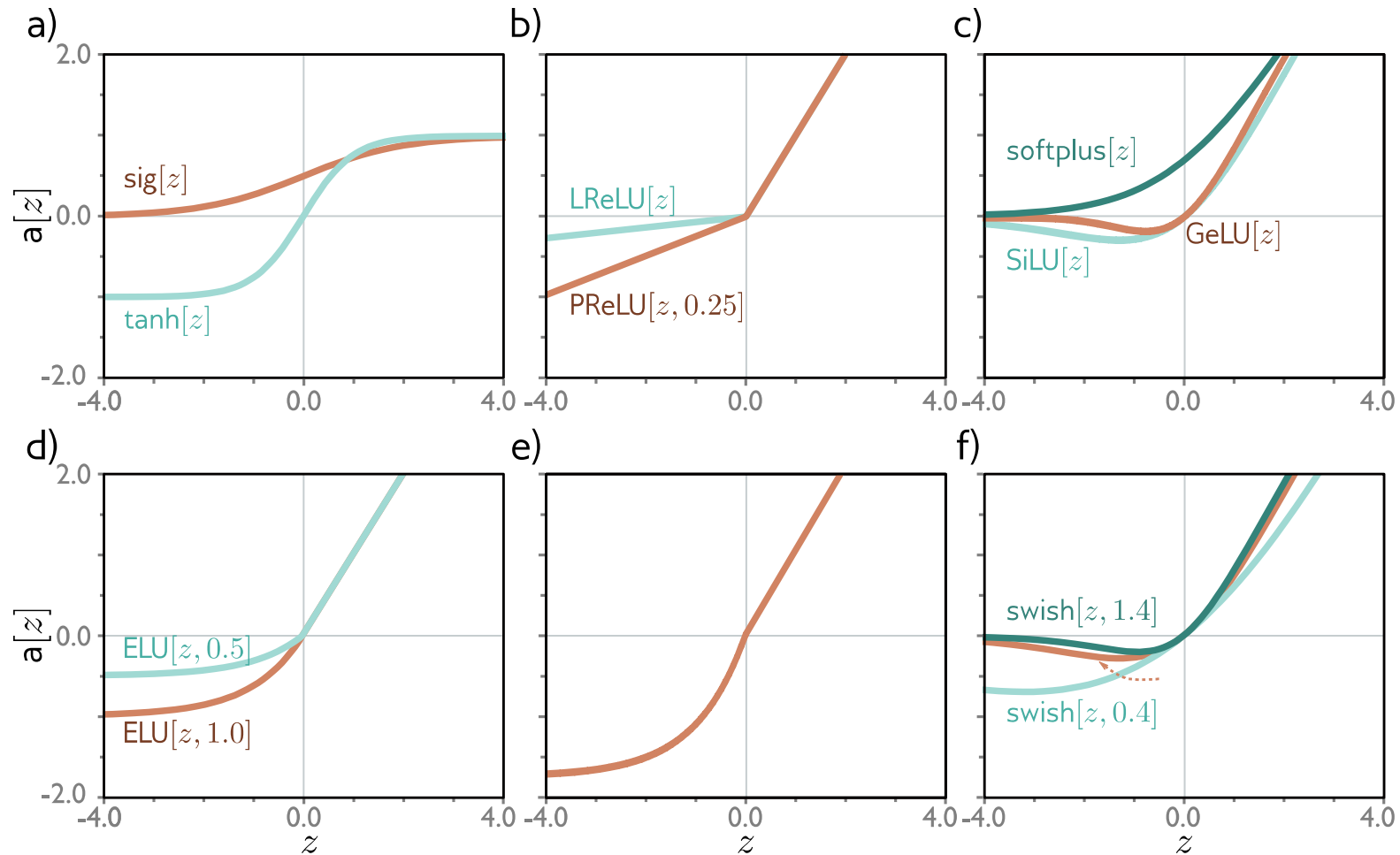- Terminology

# Nomenclature



Input layer

Hidden layer

Output layer

$x_1$

$x_2$

$x_3$

$h_1$

$h_2$

$h_3$

$h_4$

$y_1$

$y_2$

Weight or
parameter

Neuron or
hidden unit

# Nomenclature



- Y-offsets = biases
- Slopes = weights
- Everything in one layer connected to everything in the next = fully connected network (multi-layer perceptron)
- No loops = feedforward network
- Values after ReLU (activation functions) = activations
- Values before ReLU = pre-activations
- One hidden layer = shallow neural network
- More than one hidden layer = deep neural network
- Number of hidden units $\approx$ capacity

# Other activation functions



a)

sig[$z$]

tanh[$z$]

b)

LReLU[$z$]

PReLU[$z, 0.25$]

c)

softplus[$z$]

SiLU[$z$]

GeLU[$z$]

d)

ELU[$z, 0.5$]

ELU[$z, 1.0$]

e)

f)

swish[$z, 1.4$]

swish[$z, 0.4$]

Ramachandran, P.,
Zoph, B., & Le, Q. V.
(2017). Searching for
activation functions.
*arXiv:1710.05941*.

65

# Any questions?