

Shallow Neural Networks

January 8-9th, 2025

Deep Learning (CS60010)

Regression

Real world input

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

Model
input

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

Model



Supervised learning
model

Model
output

$$[340]$$

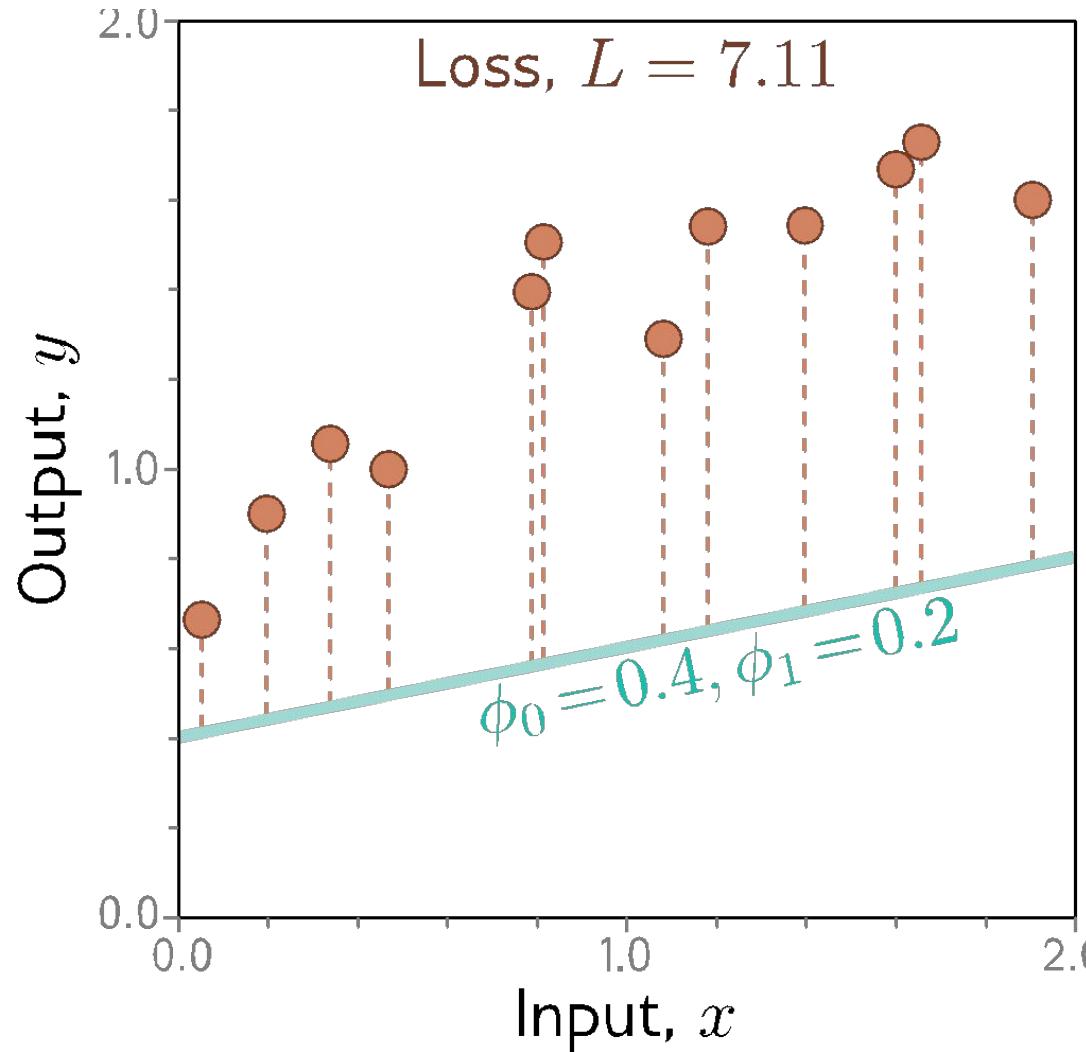
P

Real world output

Predicted price
is \$340k

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

Example: 1D Linear regression loss function

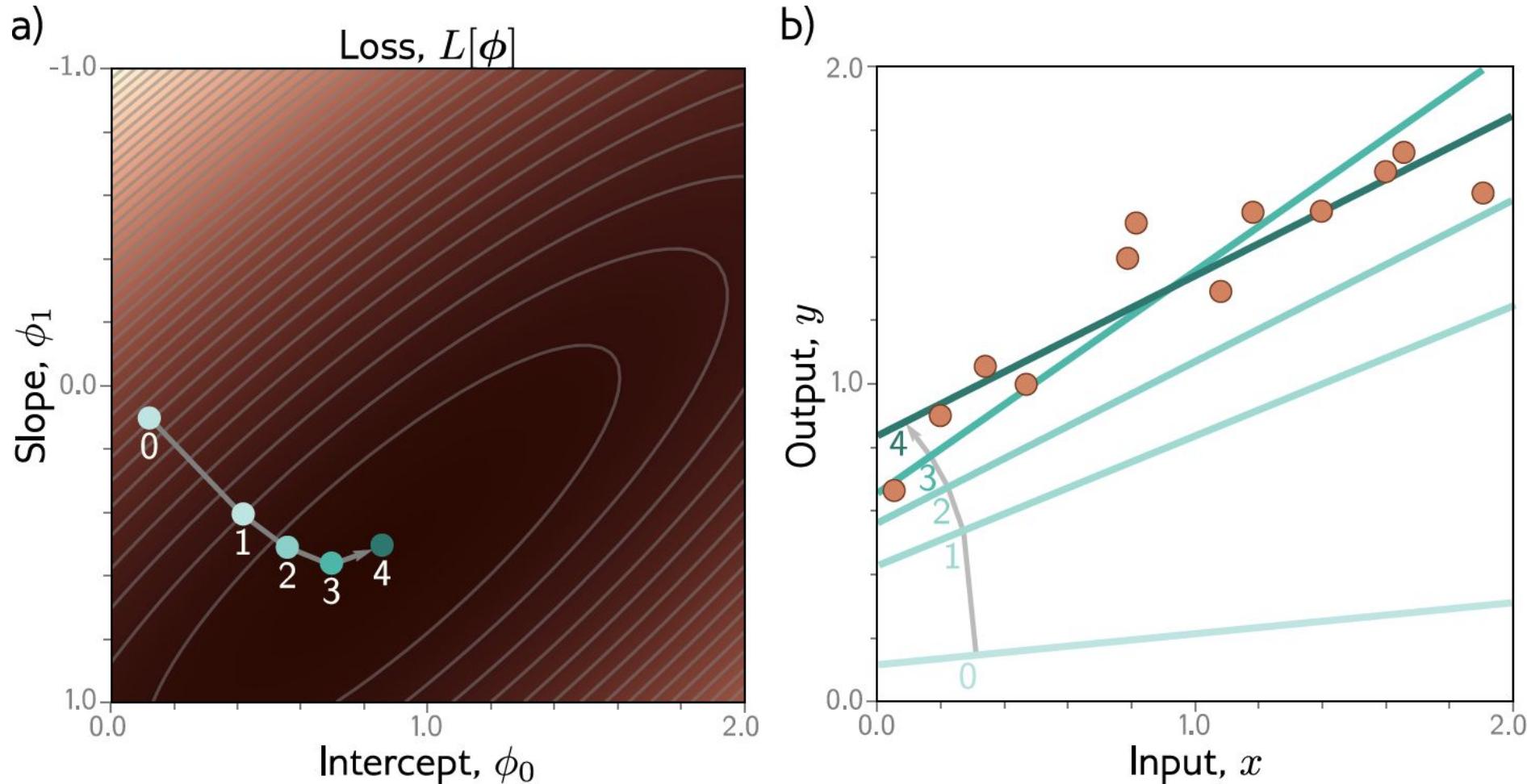


Loss function:

$$\begin{aligned} L[\phi] &= \sum_{i=1}^I (f[x_i, \phi] - y_i)^2 \\ &= \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2 \end{aligned}$$

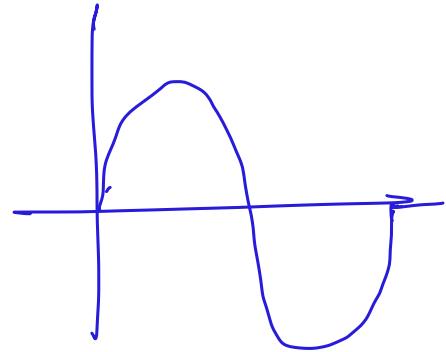

“Least squares loss function”

Example: 1D Linear regression training



This technique is known as **gradient descent**

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

$$\begin{aligned}y &= f[x, \phi] \\&= \phi_0 + \phi_1 x\end{aligned}$$

Example shallow network

$$\begin{aligned}y &= f[x, \phi] \\&= \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]\end{aligned}$$

Number of hidden layers?

Example shallow network

$$h_i = a [\theta_{i0} + \theta_{ii} x]$$

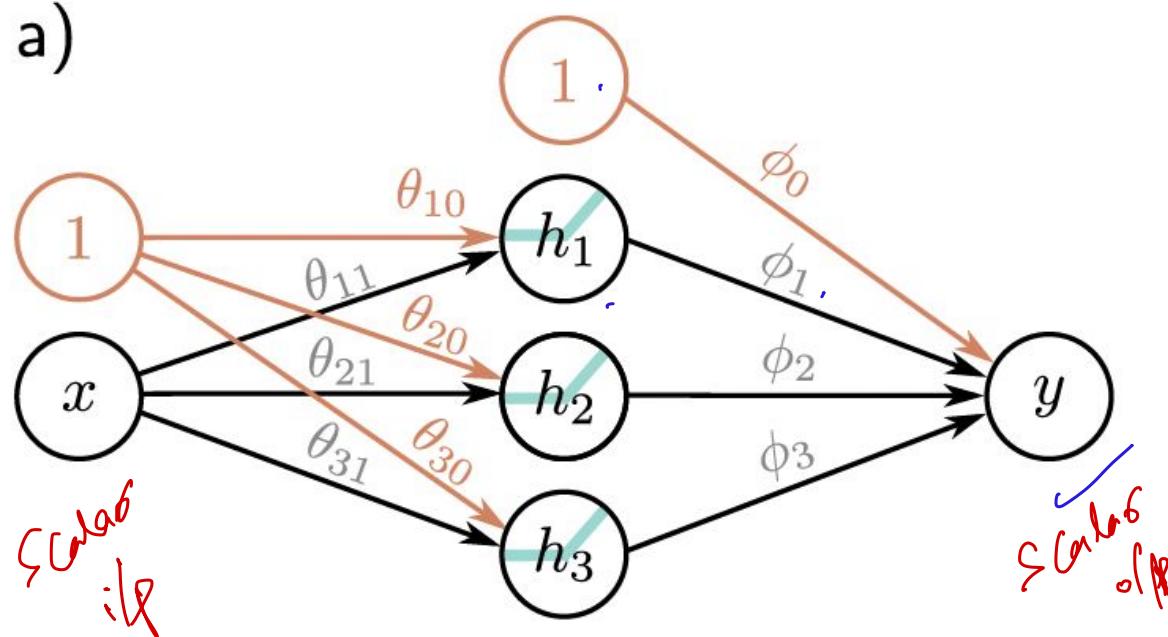
$$\phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

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

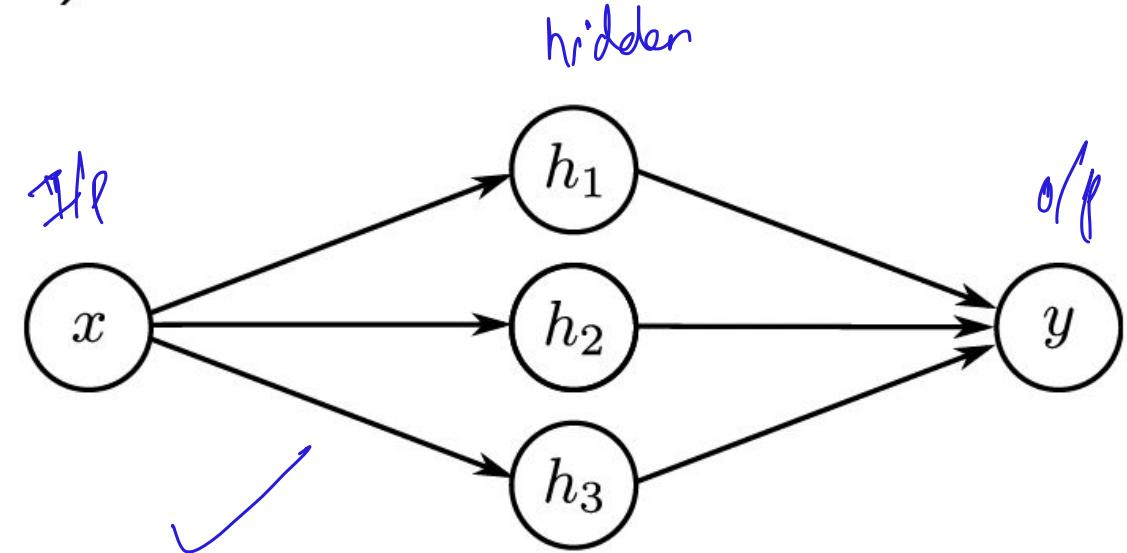
✓

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

a)



b)



Example shallow network

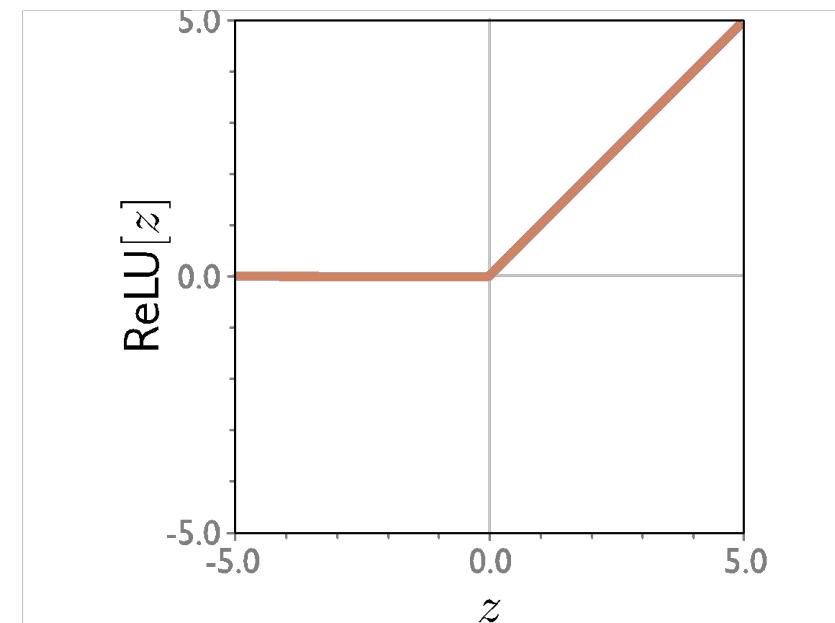
$$y = f[x, \phi]$$

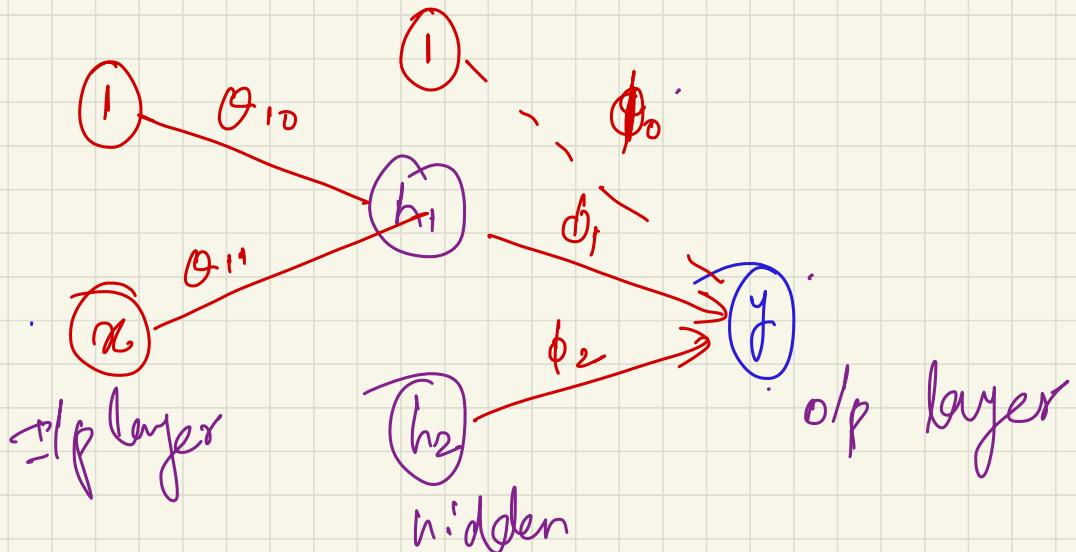
$$= \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$$

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

Rectified Linear Unit
(particular kind of activation function)

Activation function





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

$$= \phi_0 + \phi_1 (a[\theta_{10} + \theta_{11}x]) + \phi_2 (a[\theta_{20} + \theta_{21}x])$$

Example shallow network

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

Annotations: Red checkmarks are placed above each term in the equation, and red arrows point from the parameters $\phi_0, \phi_1, \phi_2, \phi_3$ to their corresponding terms.

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 particular function
 - Given parameters can perform inference (run equation)
 - Given training dataset $\{x_i, y_i\}_{i=1}^I$
 - Define loss function $L[\phi]$ (least squares)
 - Change parameters to minimize loss function
- Annotations: Red arrows point from the text "Given parameters can perform inference" to the word "inference", from "Given training dataset" to the set $\{x_i, y_i\}_{i=1}^I$, and from "Define loss function" to the symbol $L[\phi]$.

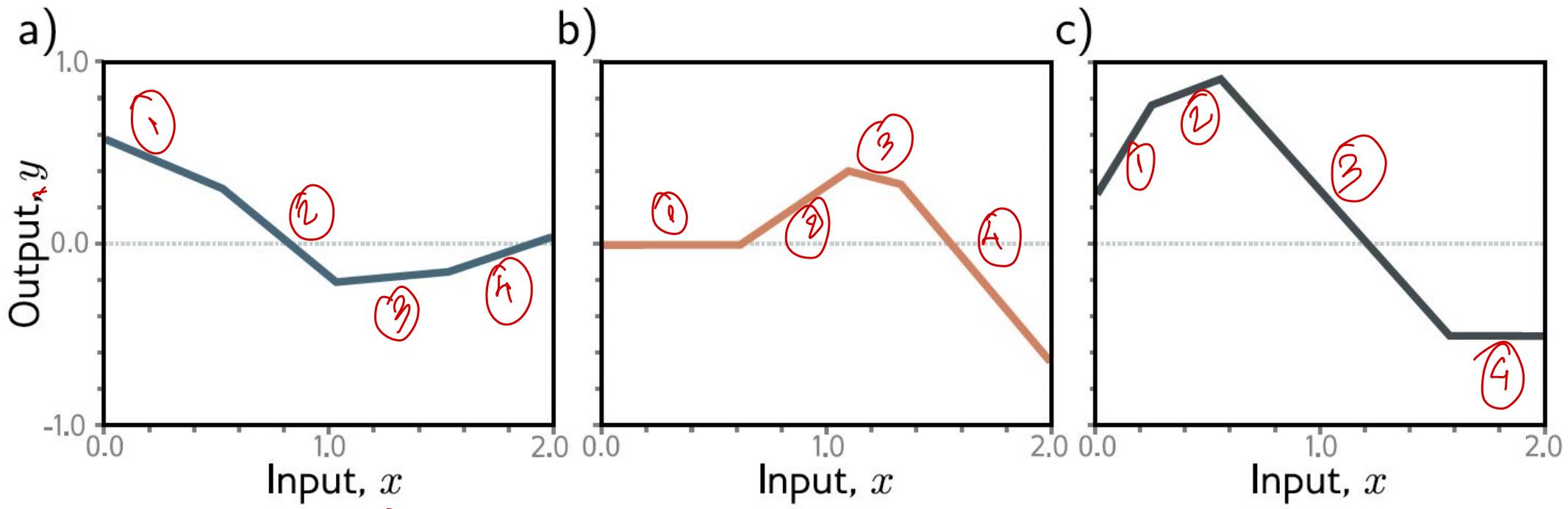
Example shallow network

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

Example shallow network

$$y = f^{(n)}$$

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



Family of functions defined by the equation for three choices of parameters

Hidden units

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

$\underbrace{\hspace{10em}}$
 h_1

$\underbrace{\hspace{10em}}$
 h_2

$\underbrace{\hspace{10em}}$
 h_3

Break down into two parts:

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

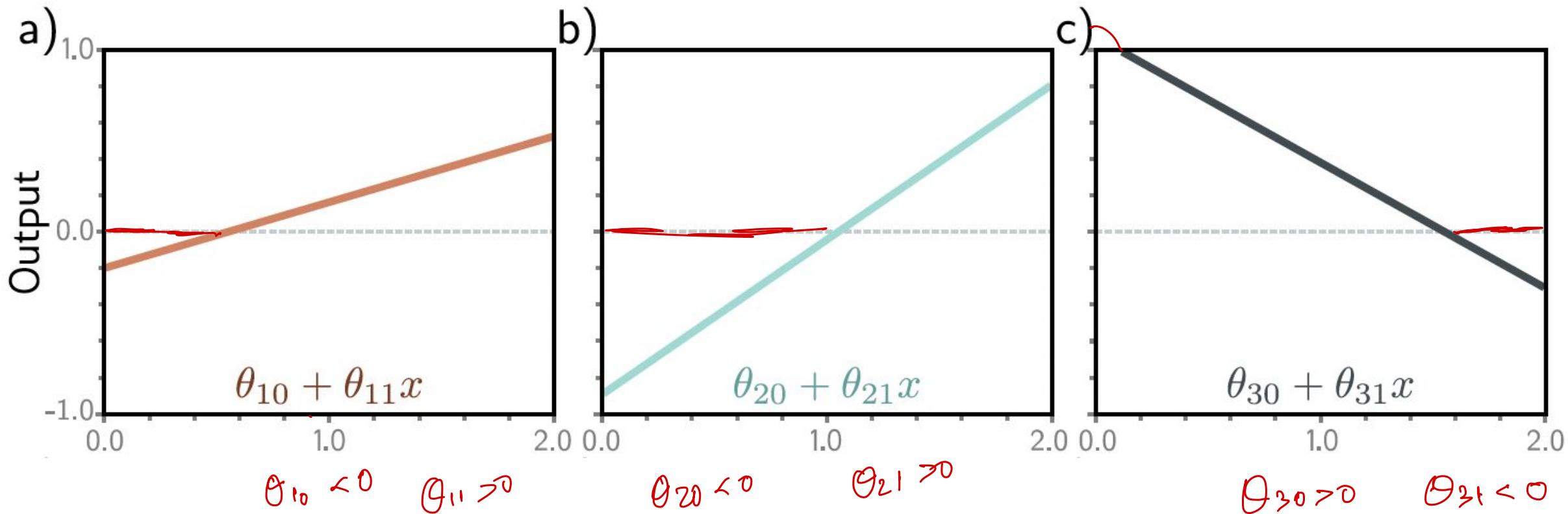
where:

Hidden units

$$\left[\begin{array}{lcl} h_1 & = & a[\theta_{10} + \theta_{11}x] \\ h_2 & = & a[\theta_{20} + \theta_{21}x] \\ h_3 & = & a[\theta_{30} + \theta_{31}x] \end{array} \right]$$

1. compute three linear functions

pre-activation

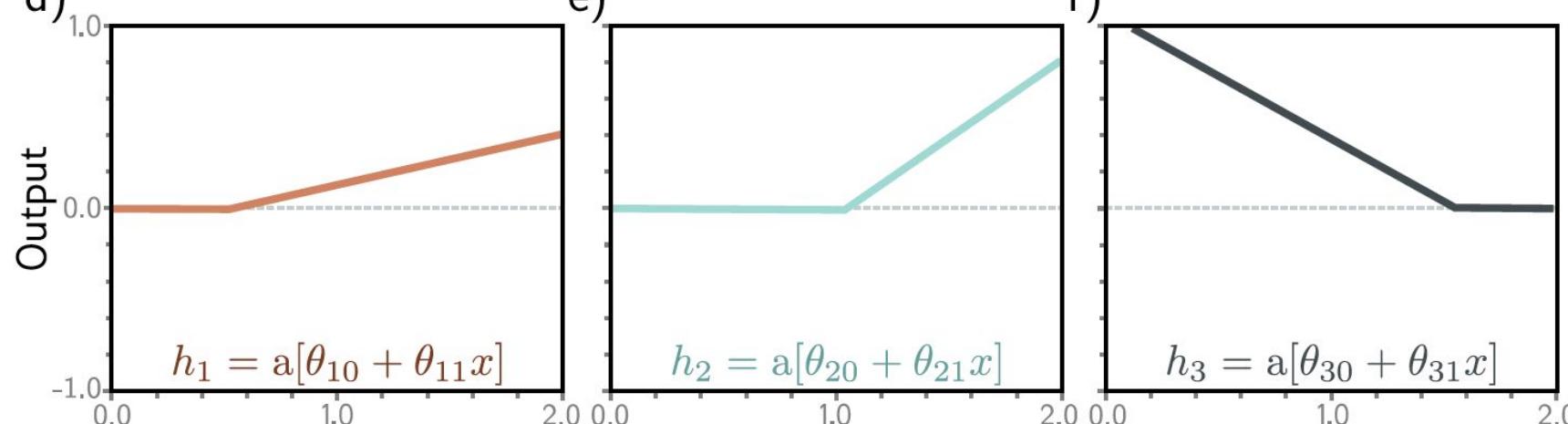
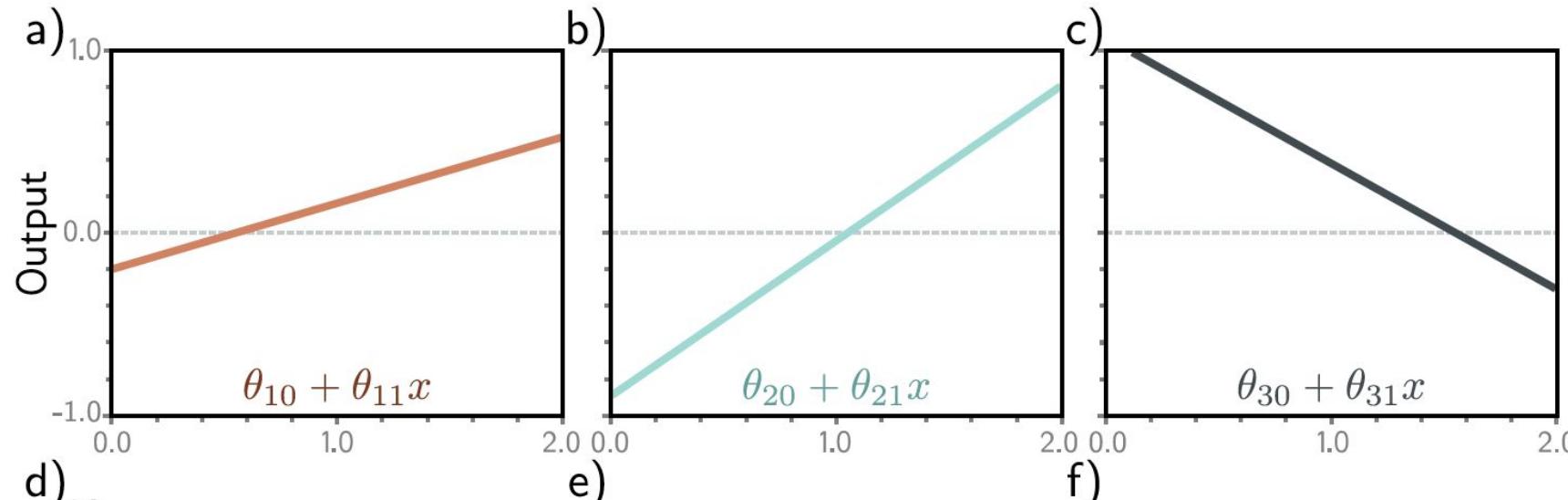


2. Pass through ReLU functions (creates hidden units)

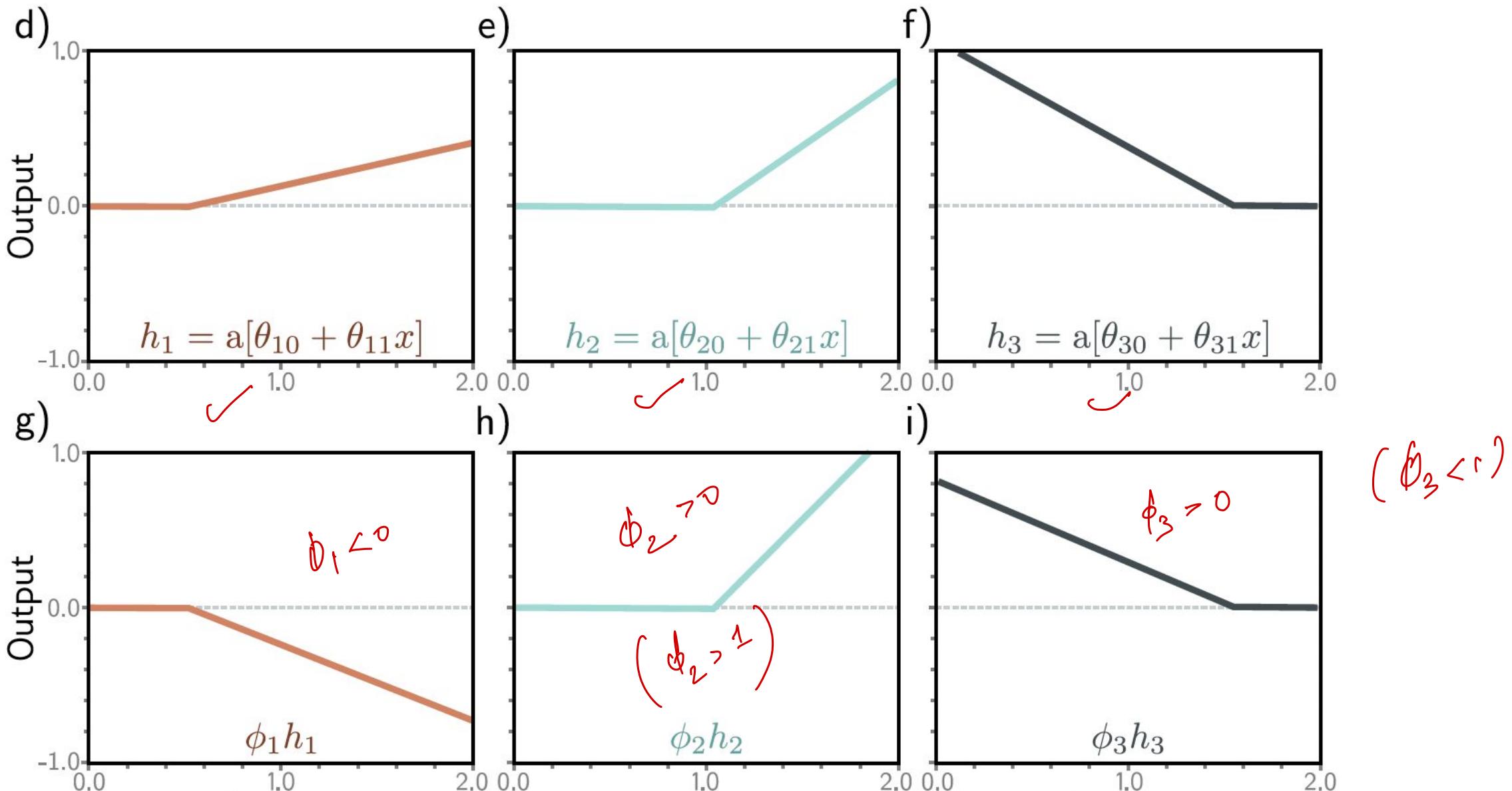
$$h_1 = a[\theta_{10} + \theta_{11}x]$$

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

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

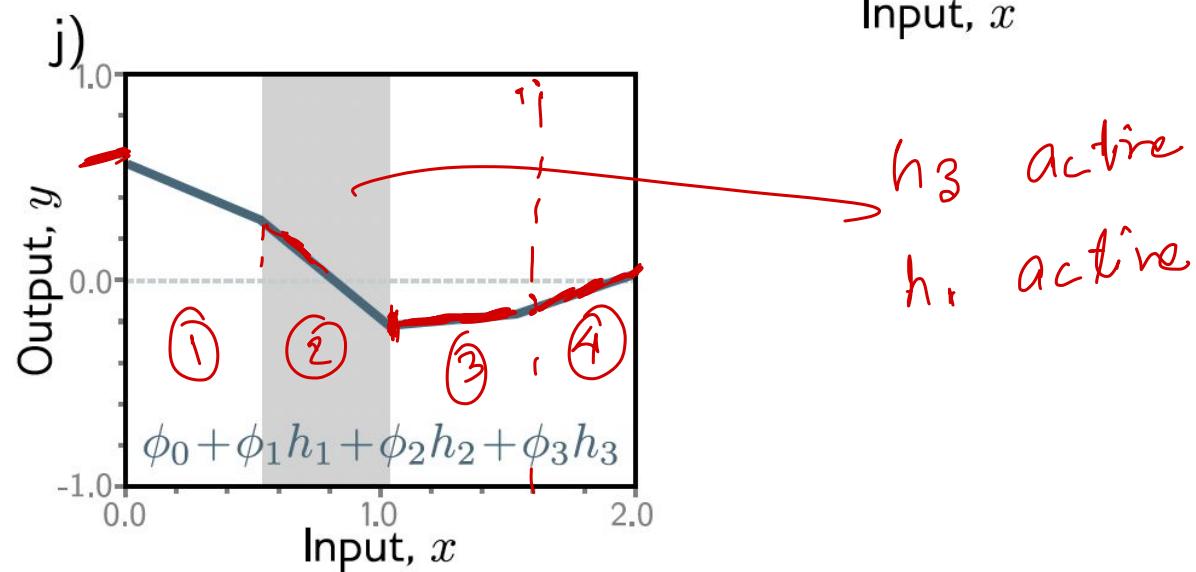
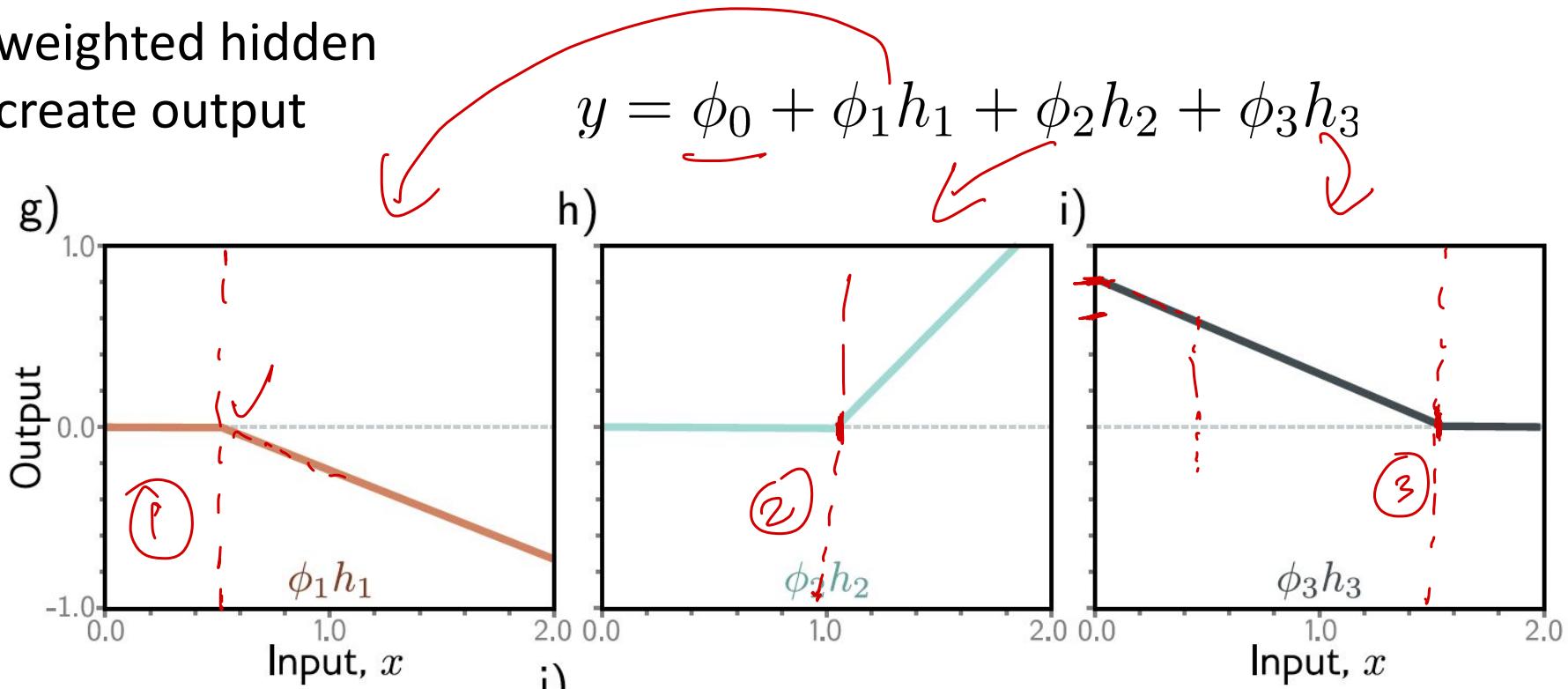


2. Weight the hidden units



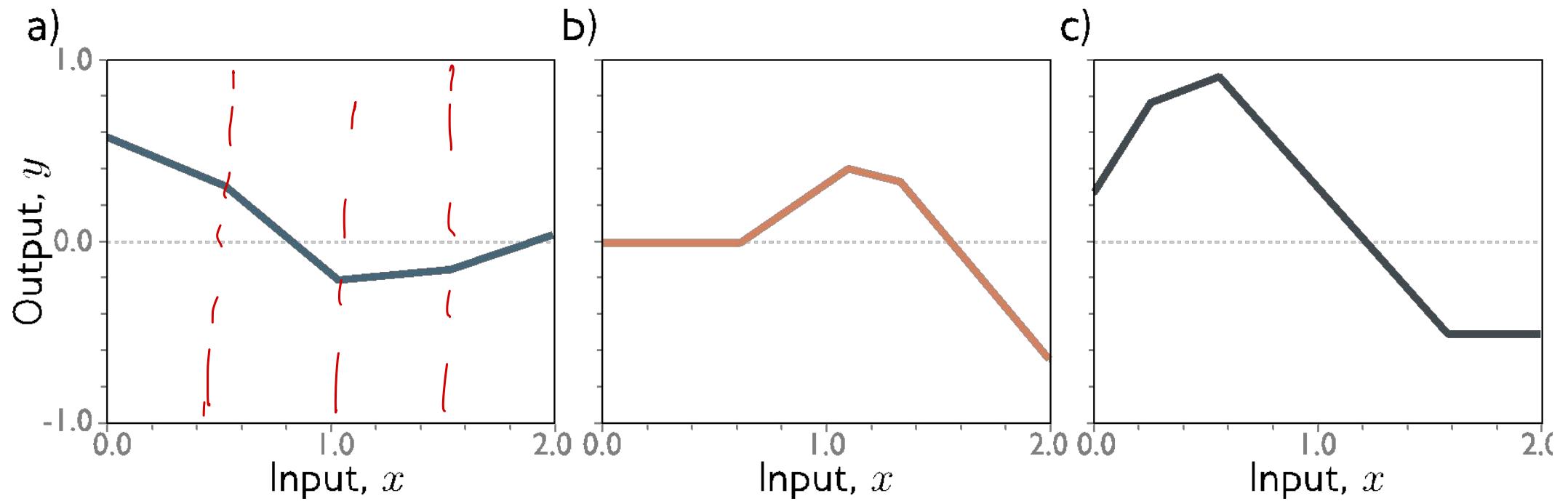
4. Sum the weighted hidden units to create output

3 joints



Example shallow network

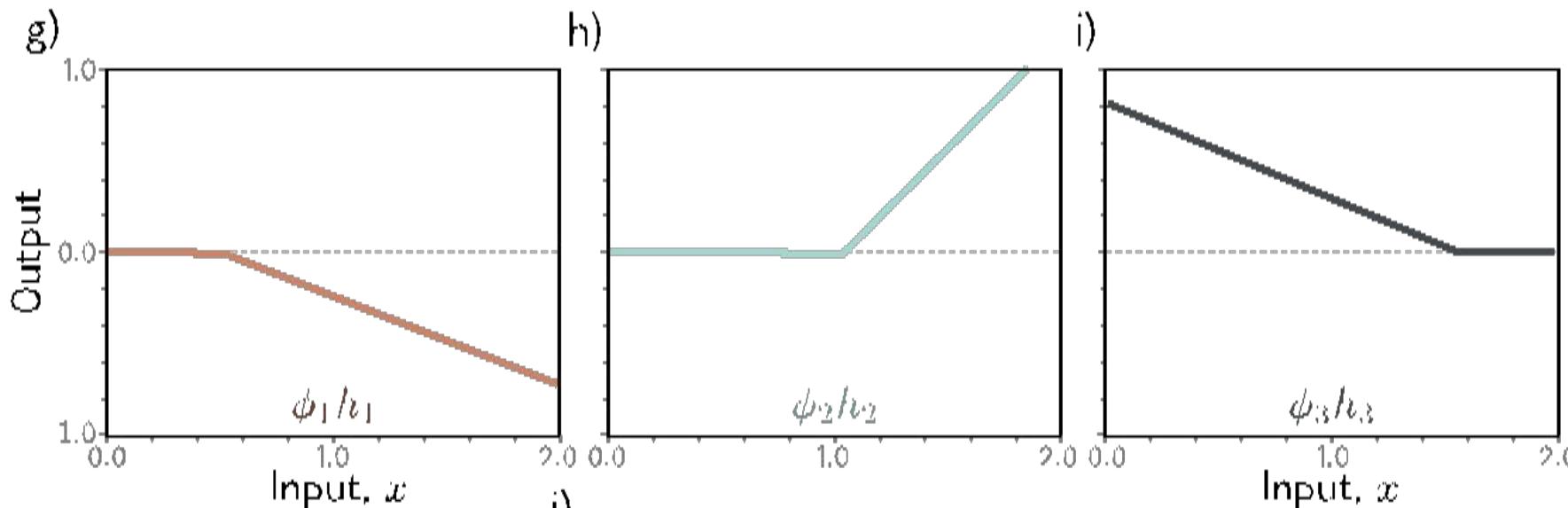
$$y = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x].$$



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

3 joints

Activation pattern = which hidden units are activated



Shaded region:

- Unit 1 active
- Unit 2 inactive
- Unit 3 active

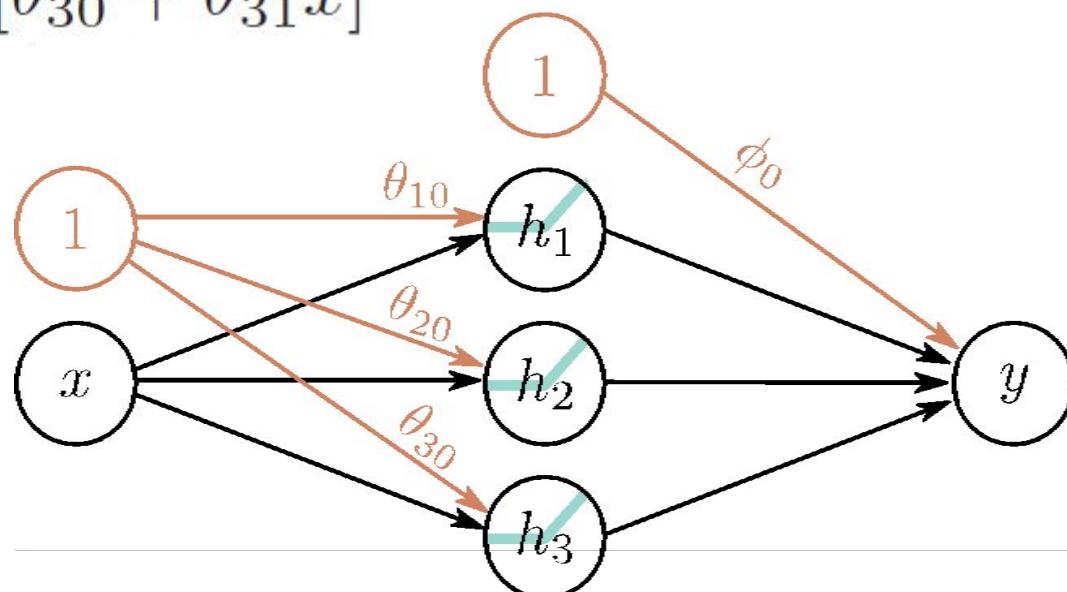
Depicting neural networks

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

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

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

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



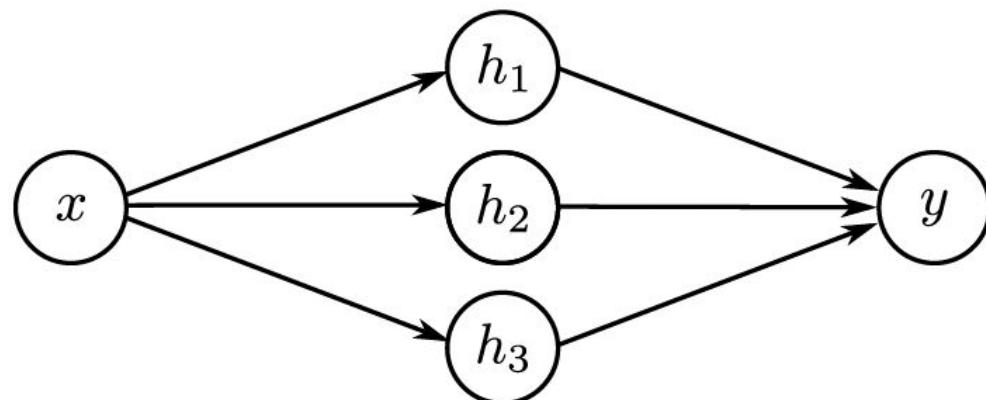
Each parameter multiplies its source and adds to its target

Depicting neural networks

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

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

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



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:

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

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

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

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

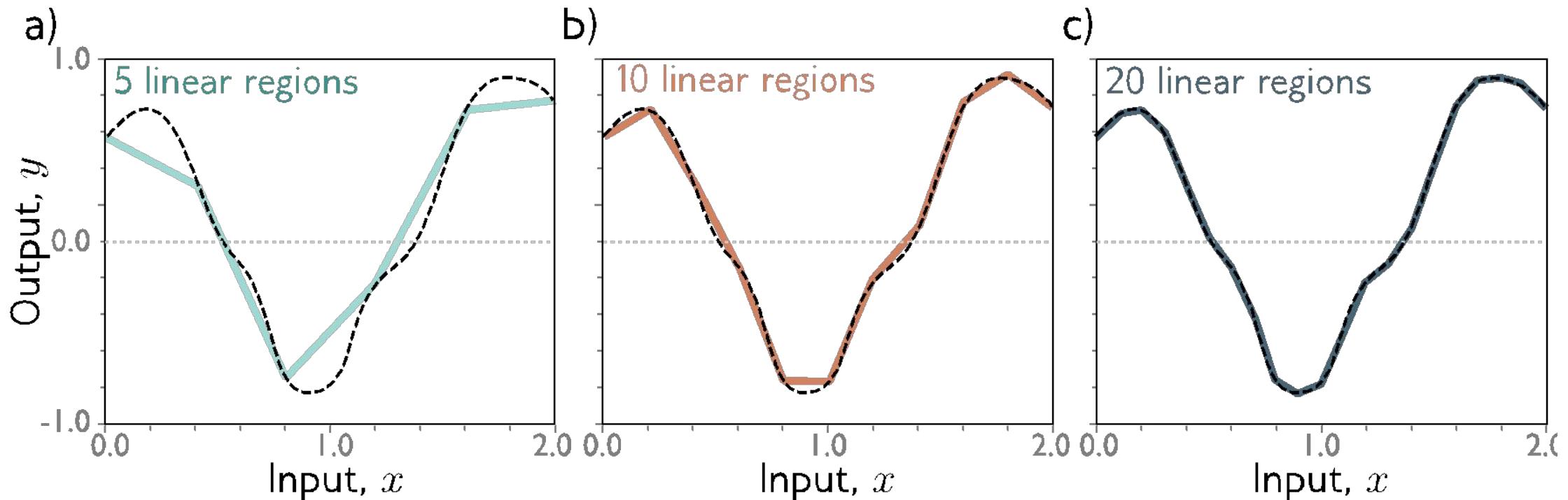
With D hidden units:

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

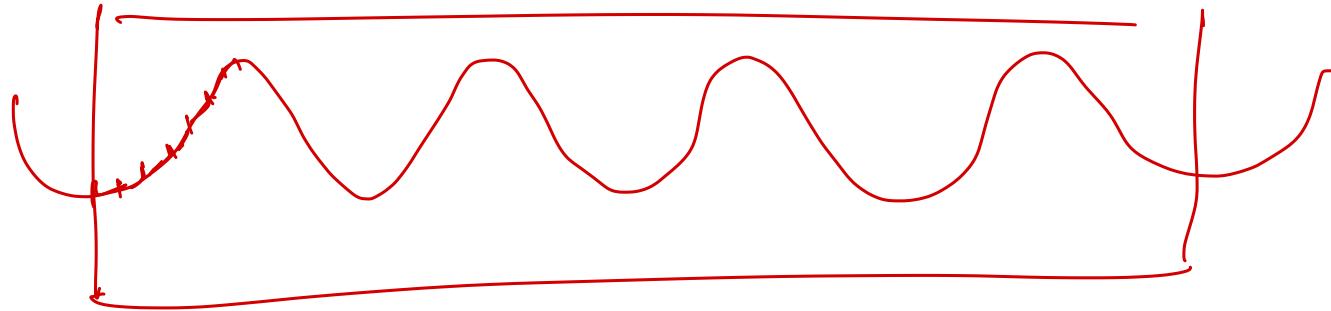
$$y = \phi_0 + \sum_{d=1}^D \phi_d h_d$$

With enough hidden units...

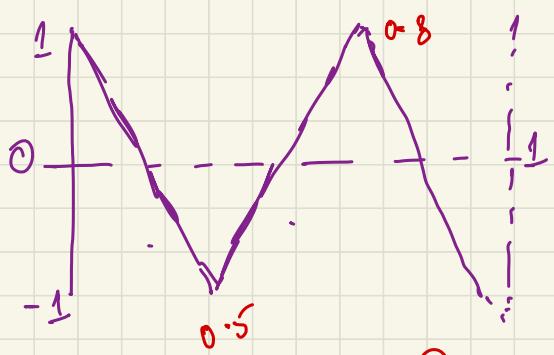
... we can describe any 1D function to arbitrary accuracy



Universal approximation theorem



“a formal proof that, with enough hidden units, a shallow neural network can describe any continuous function on a compact subset of \mathbb{R}^D to arbitrary precision”



[hw]

Can you model this using
a shallow NN with 2
neurons? If yes, give the
construction. If no, provide the proof.

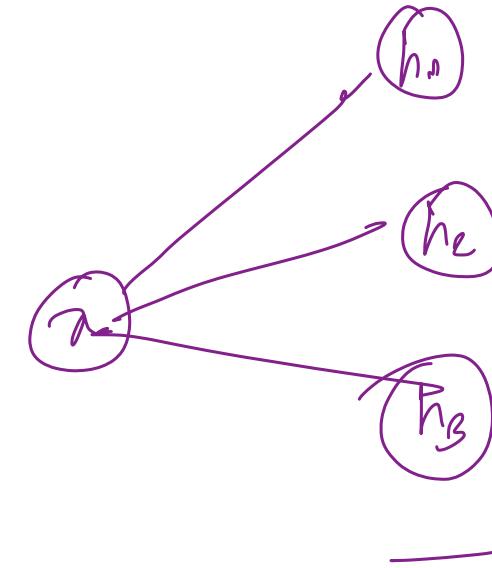
$$\phi_0 + \phi_1 h_1 + \phi_2 h_2$$



Can you model this using
3 neurons? If yes, provide
construction. If no, prove.

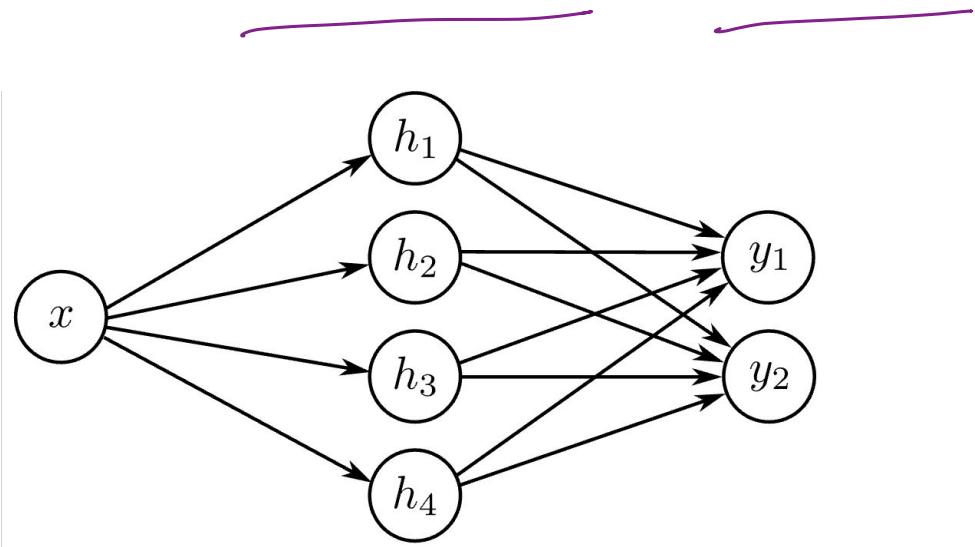
Shallow neural networks

- 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



Two outputs

- 1 input, 4 hidden units, 2 outputs

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

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

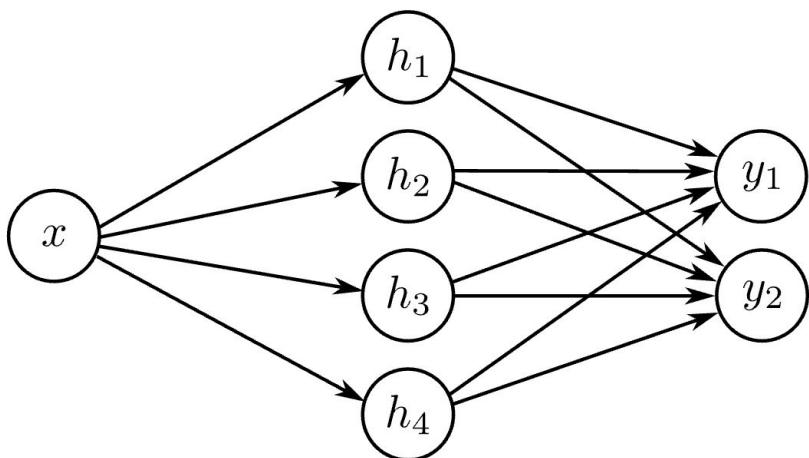
$$h_3 = a[\theta_{30} + \theta_{31}x]$$

$$h_4 = a[\theta_{40} + \theta_{41}x]$$

$$\underline{y_1} = \phi_{10} + \phi_{11}h_1 + \phi_{12}h_2 + \phi_{13}h_3 + \phi_{14}h_4$$

$$y_2 = \phi_{20} + \phi_{21}h_1 + \phi_{22}h_2 + \phi_{23}h_3 + \phi_{24}h_4$$

— - - - —



How would the outputs look like?

Two outputs

- 1 input, 4 hidden units, 2 outputs

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

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

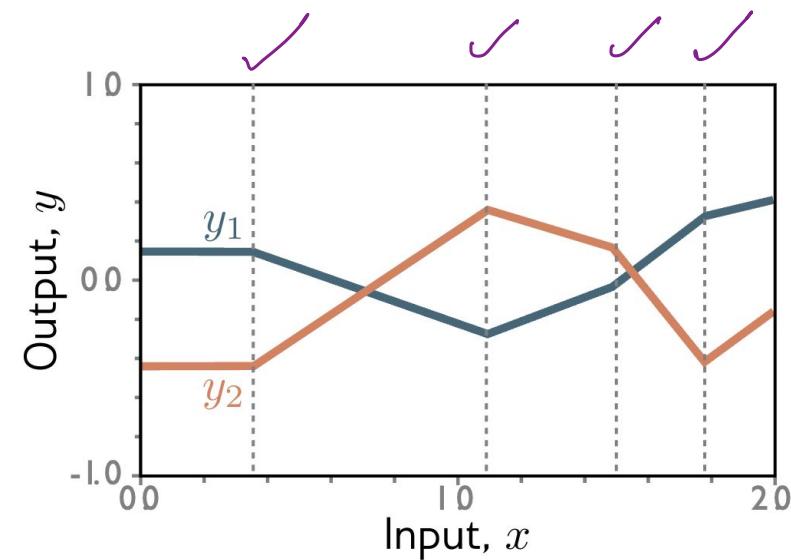
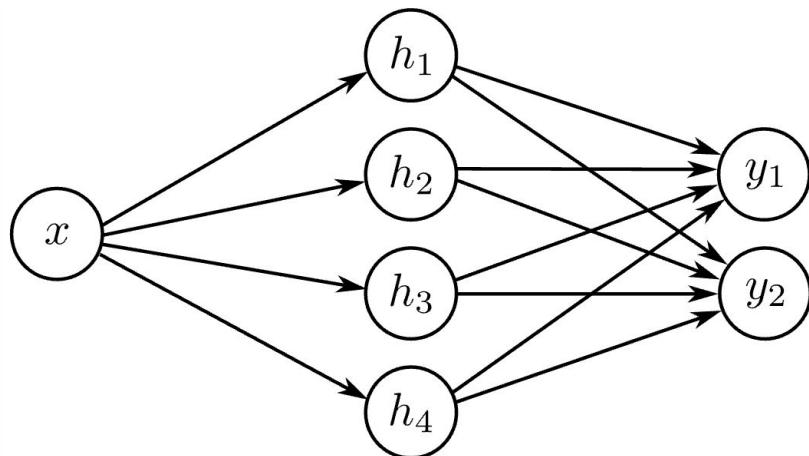
$$h_3 = a[\theta_{30} + \theta_{31}x]$$

$$h_4 = a[\theta_{40} + \theta_{41}x]$$

$$y_1 = \underline{\phi_{10}} + \phi_{11}h_1 + \phi_{12}h_2 + \phi_{13}h_3 + \phi_{14}h_4$$

$$y_2 = \underline{\phi_{20}} + \phi_{21}h_1 + \phi_{22}h_2 + \phi_{23}h_3 + \phi_{24}h_4$$

↗
0

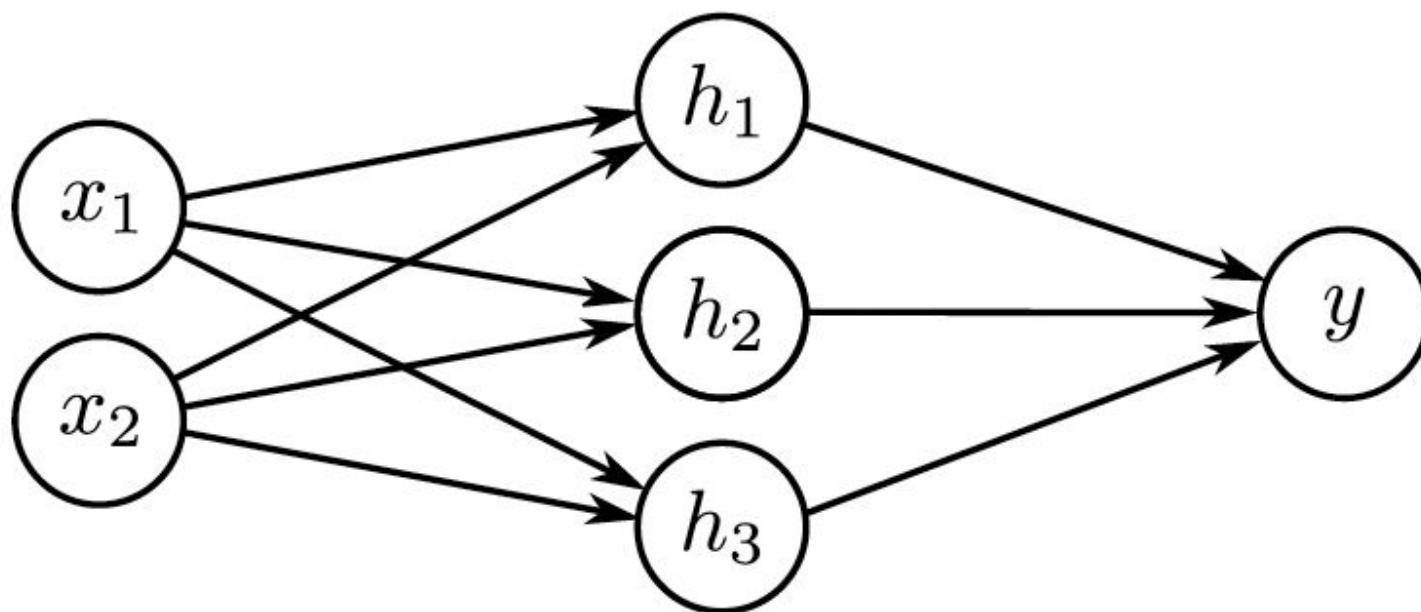


Shallow neural networks

- 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

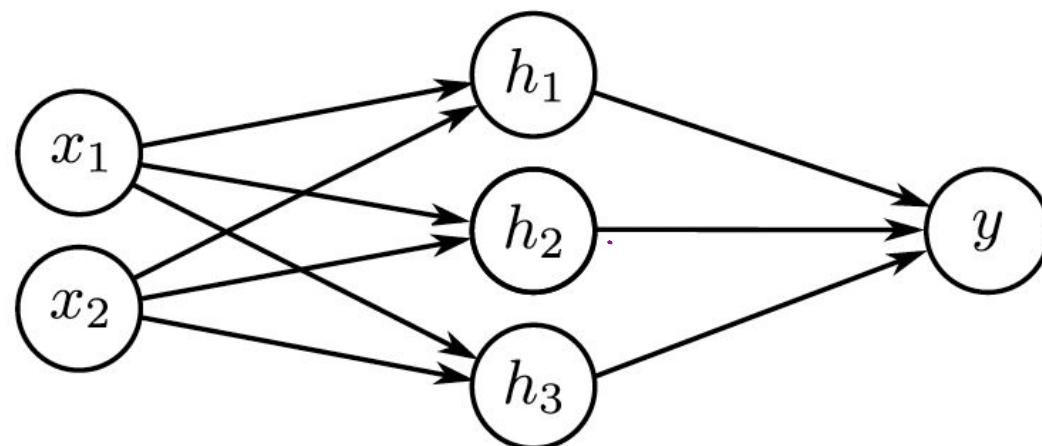


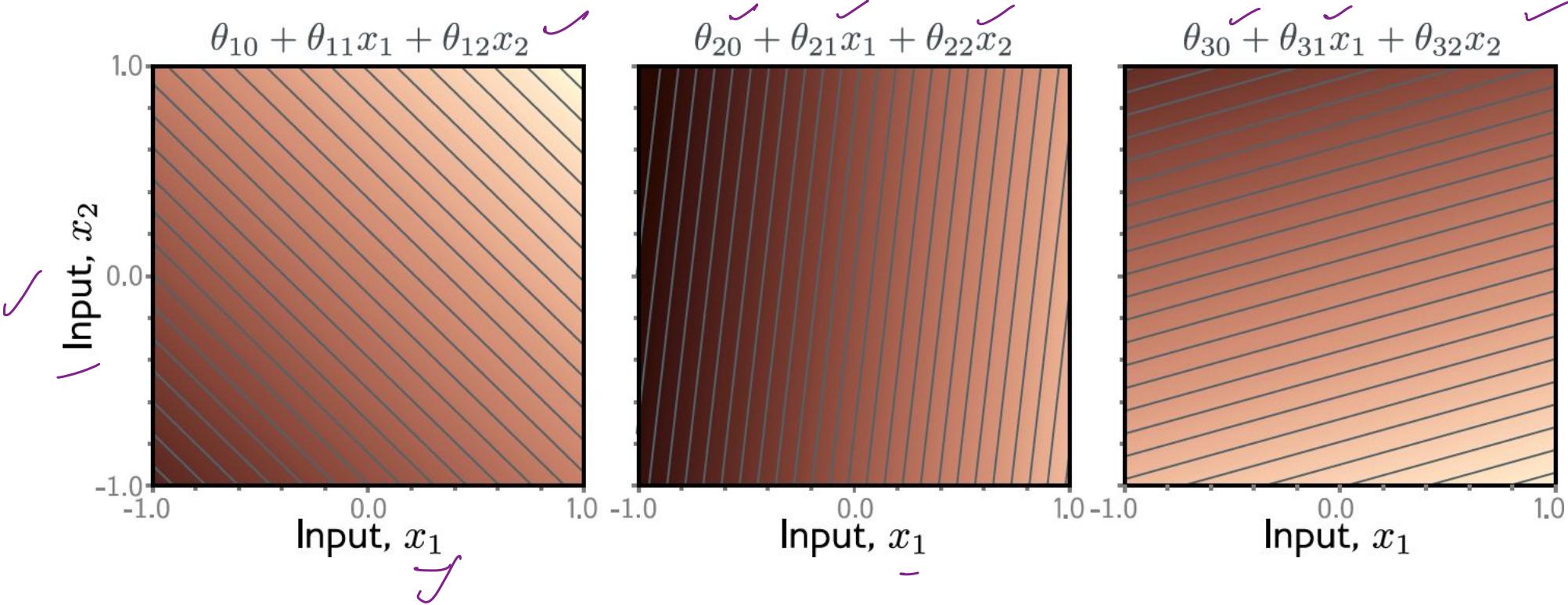
Two inputs

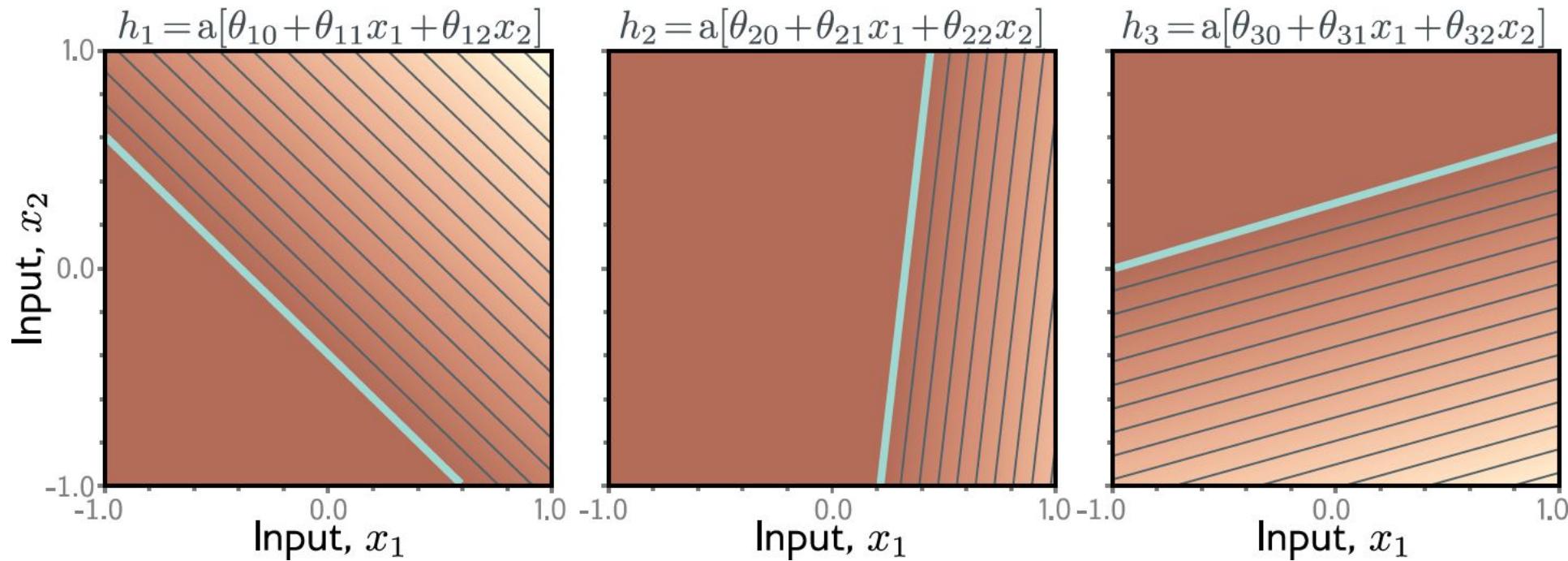
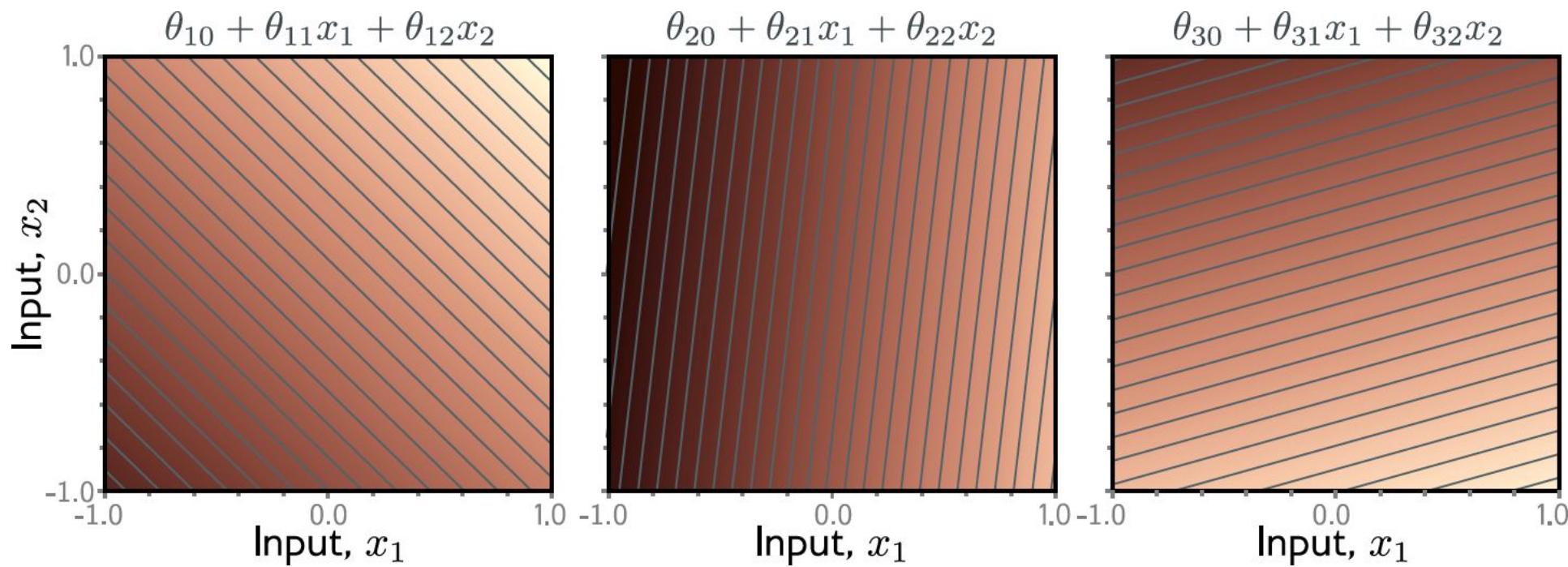
- 2 inputs, 3 hidden units, 1 output

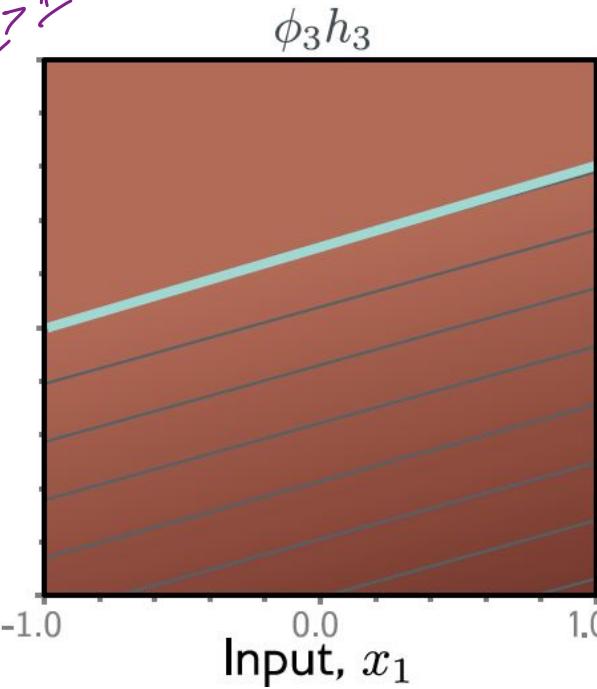
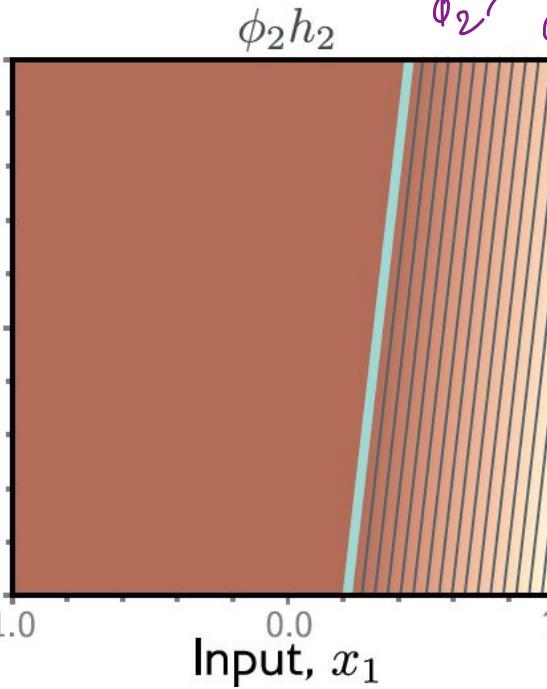
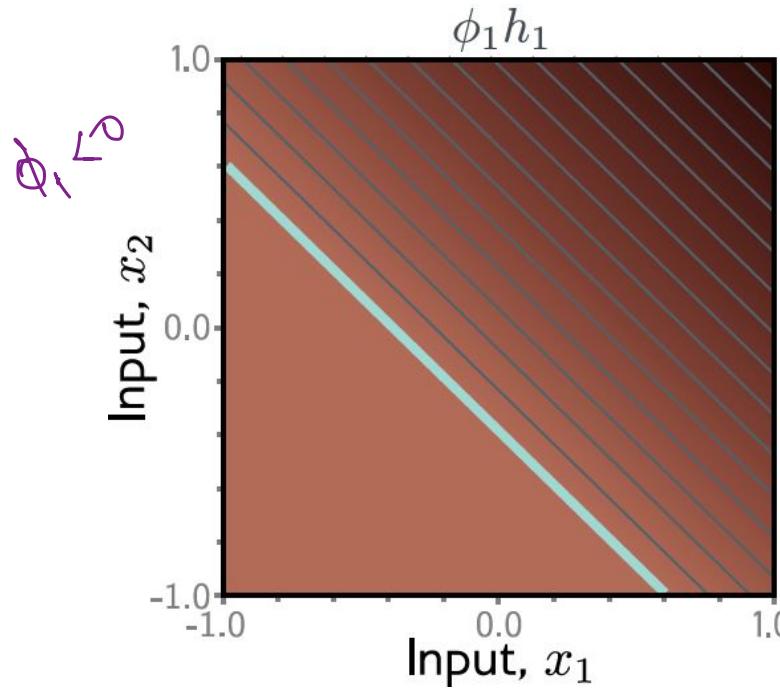
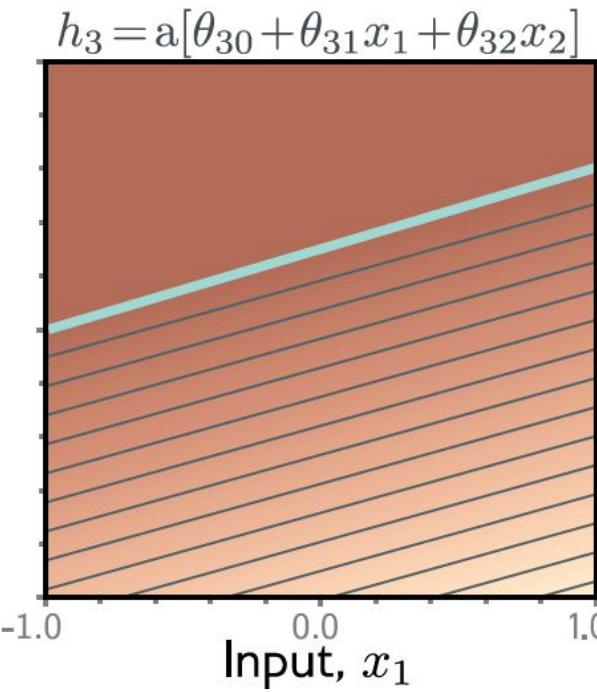
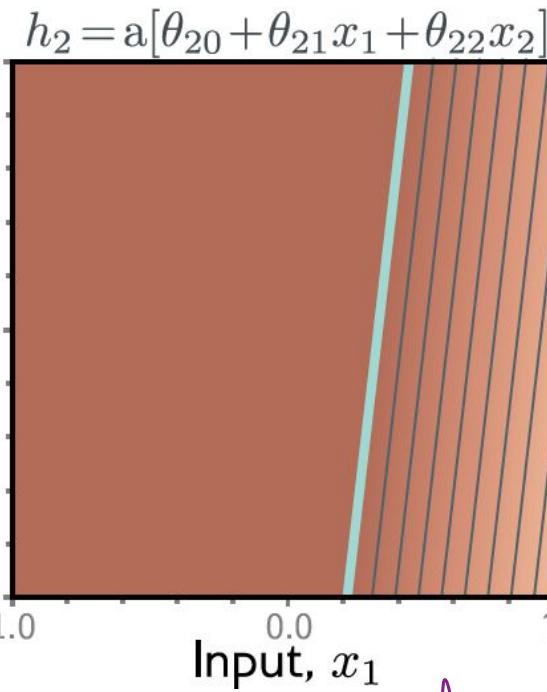
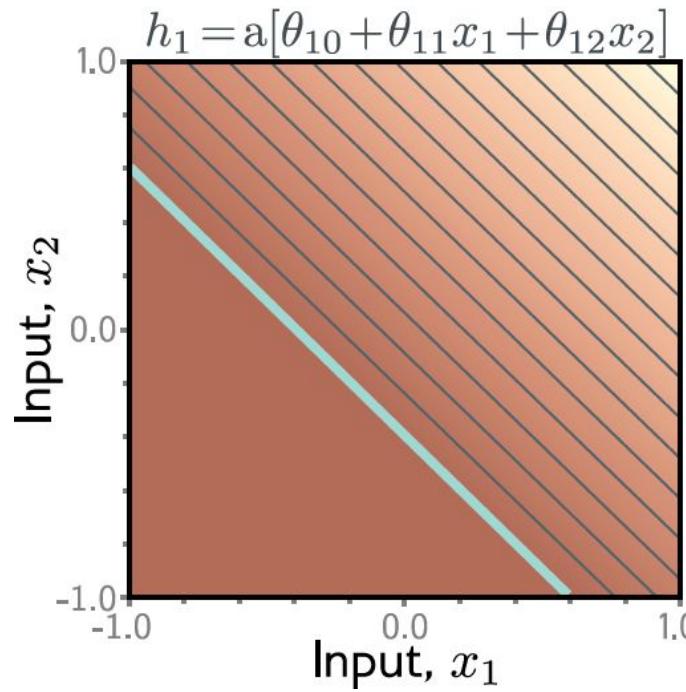
$$\begin{aligned} h_1 &= a[\theta_{10} + \theta_{11}\cancel{x_1} + \theta_{12}\cancel{x_2}] \\ h_2 &= a[\cancel{\theta_{20}} + \theta_{21}\cancel{x_1} + \theta_{22}\cancel{x_2}] \\ h_3 &= a[\cancel{\theta_{30}} + \theta_{31}x_1 + \theta_{32}x_2] \end{aligned}$$

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







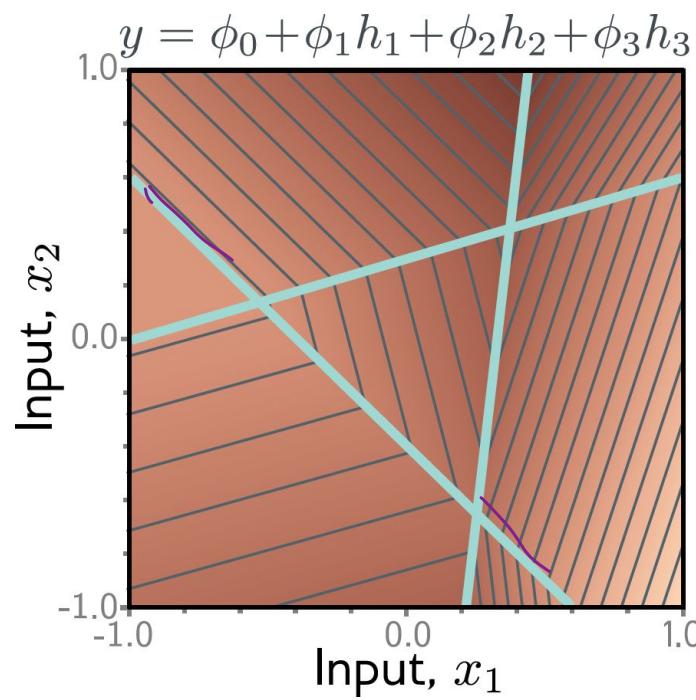
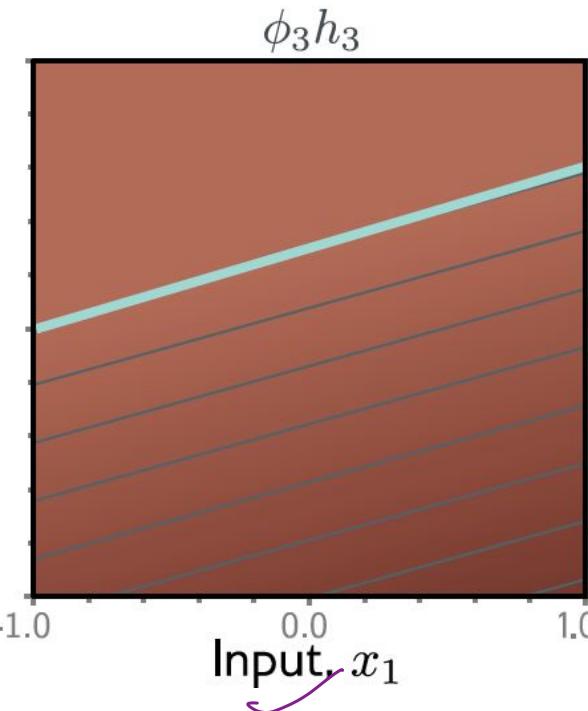
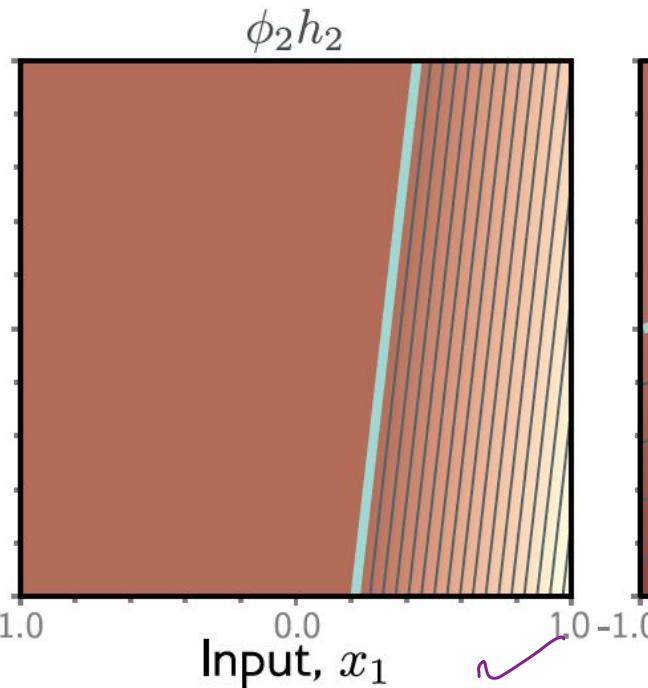
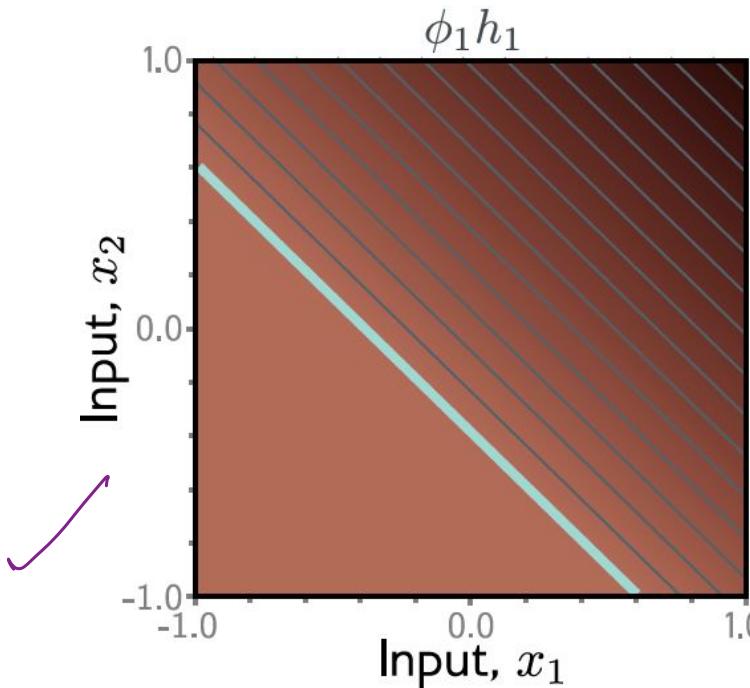


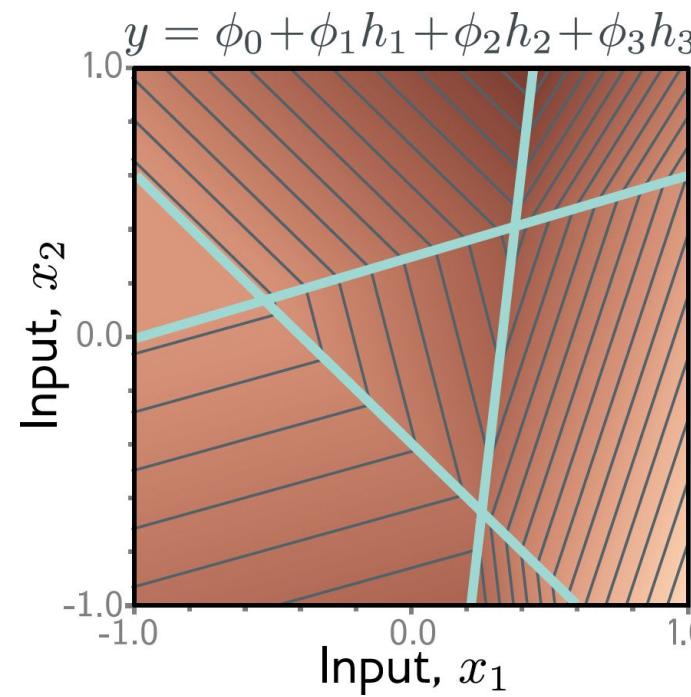
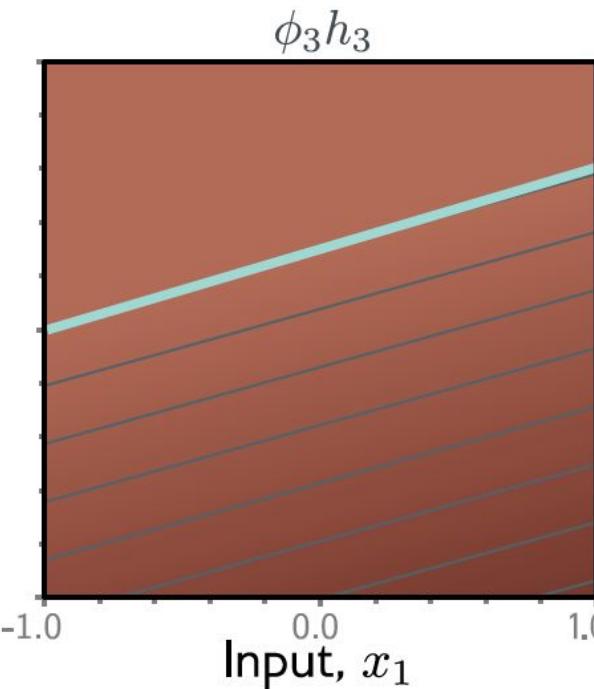
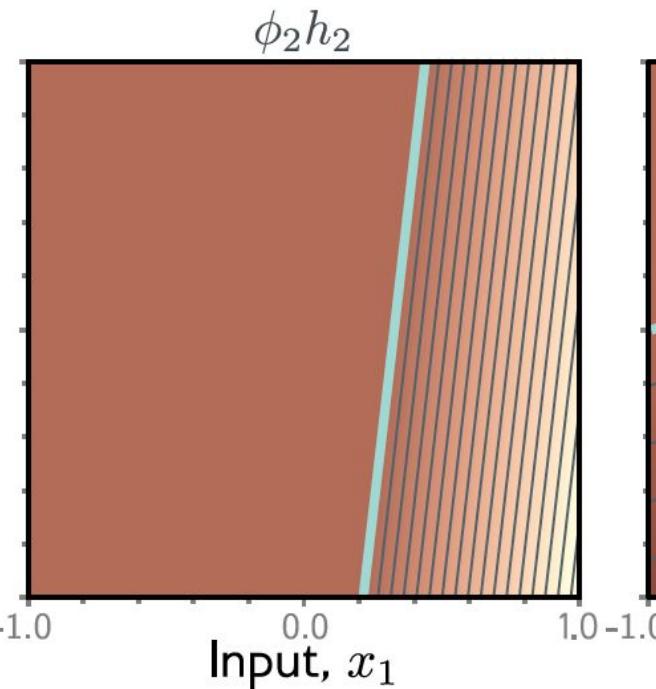
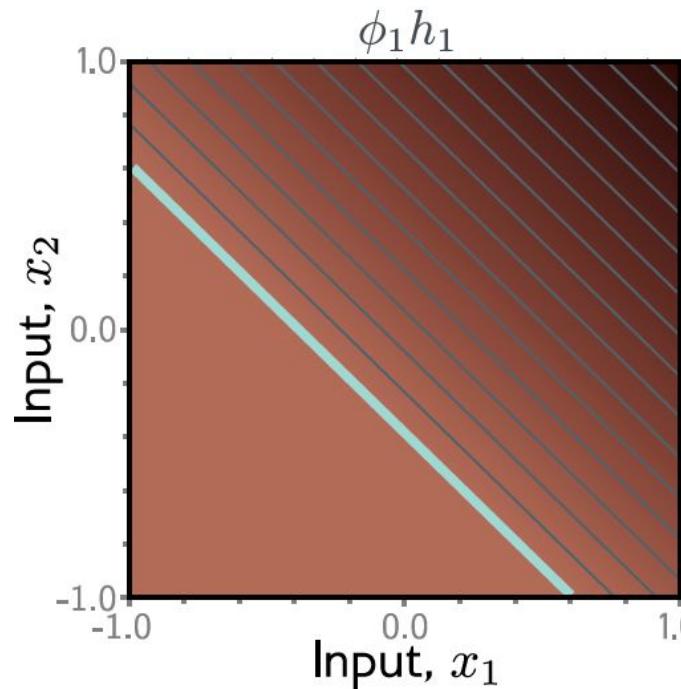
$\phi_1 < 0$

$\phi_2 > 0$

$\phi_3 > 0$

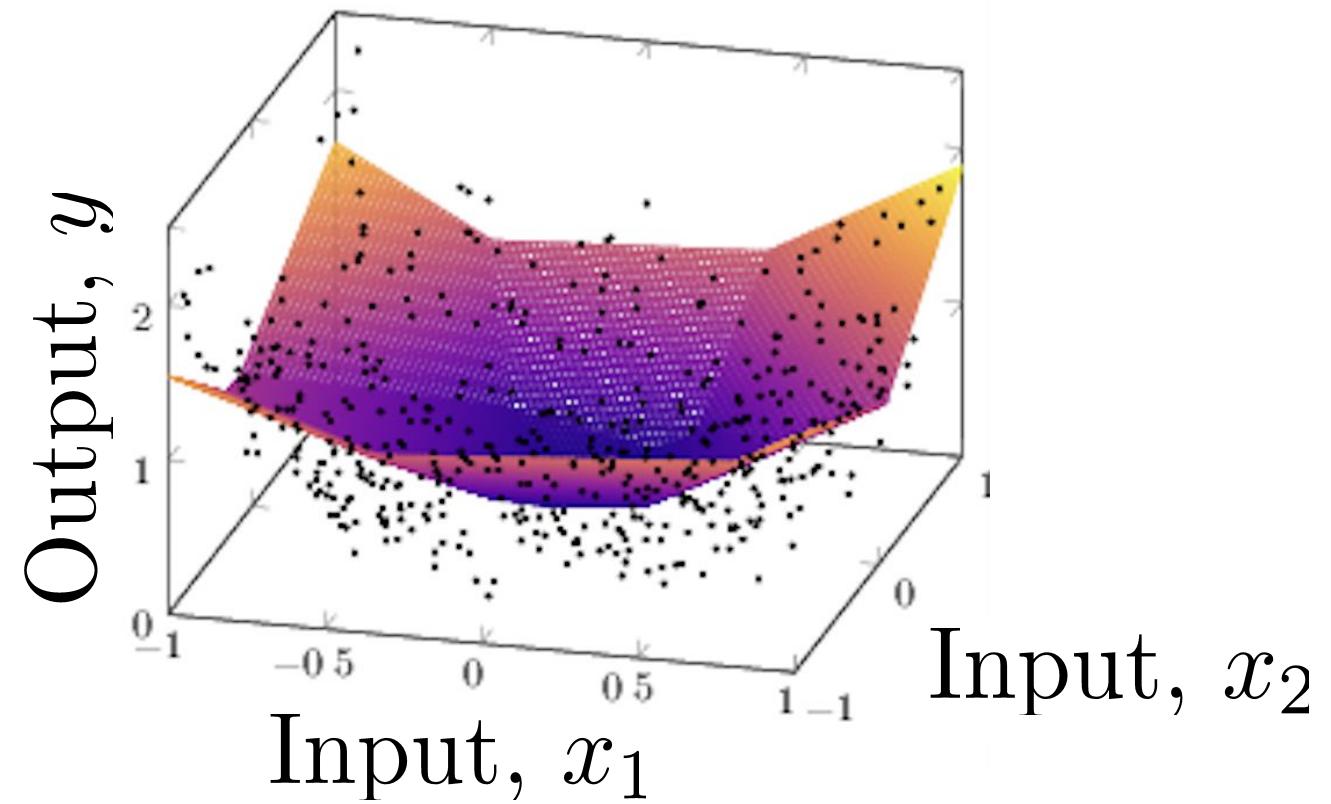
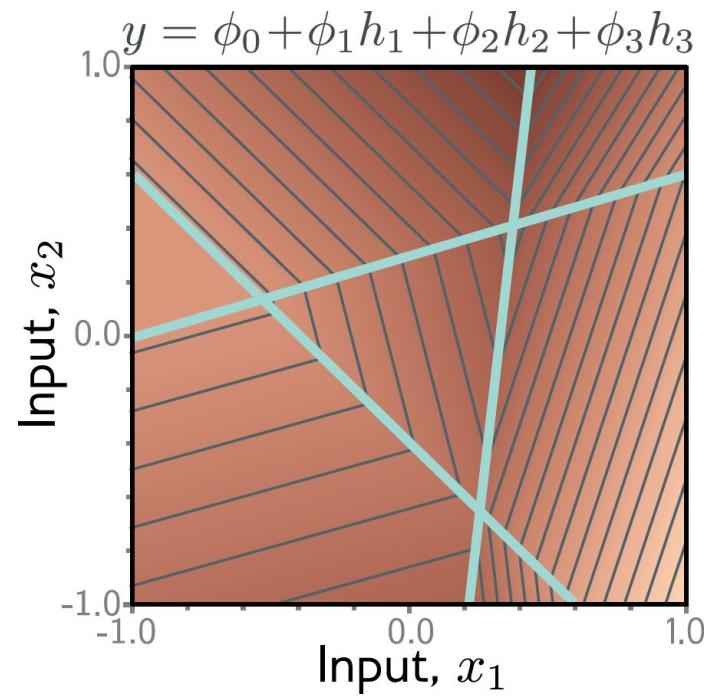
$-1 < \phi_3 < 0$





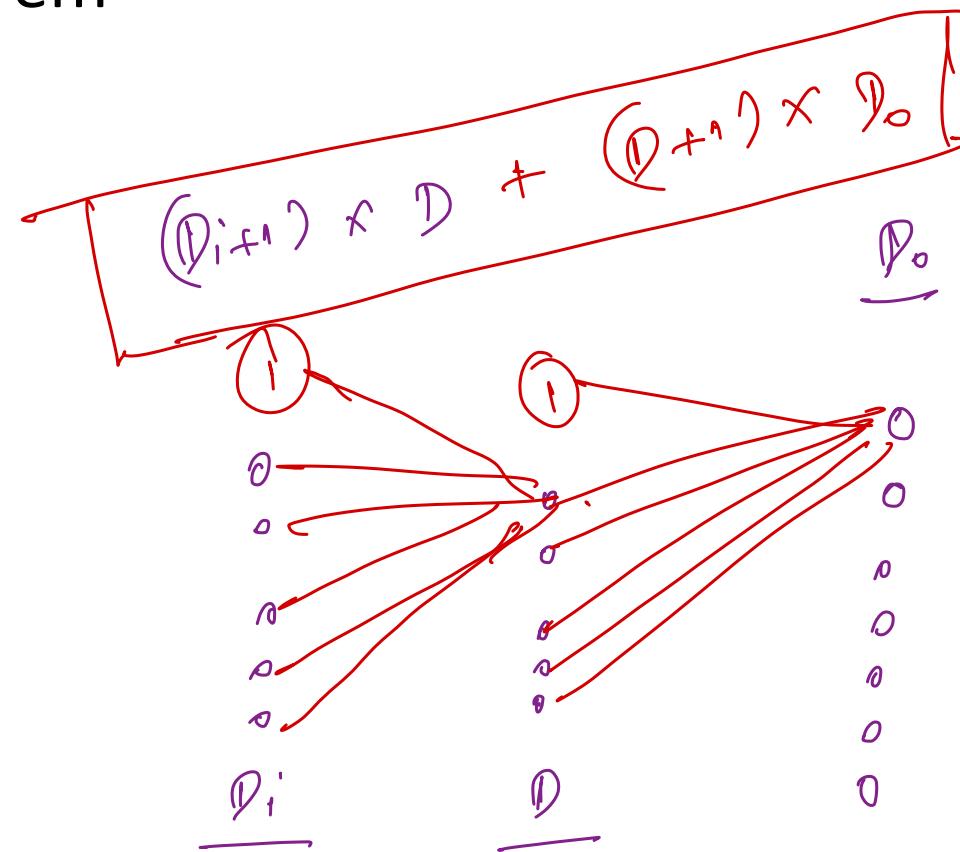
Convex polygons

Fitting



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

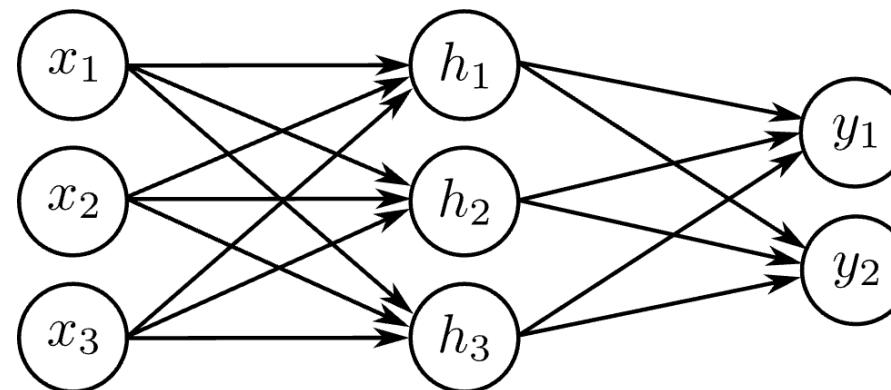


Arbitrary inputs, hidden units, outputs

- $\underbrace{D_o \text{ Outputs}}$, $\underbrace{D \text{ hidden units}}$, and $\underbrace{D_i \text{ inputs}}$

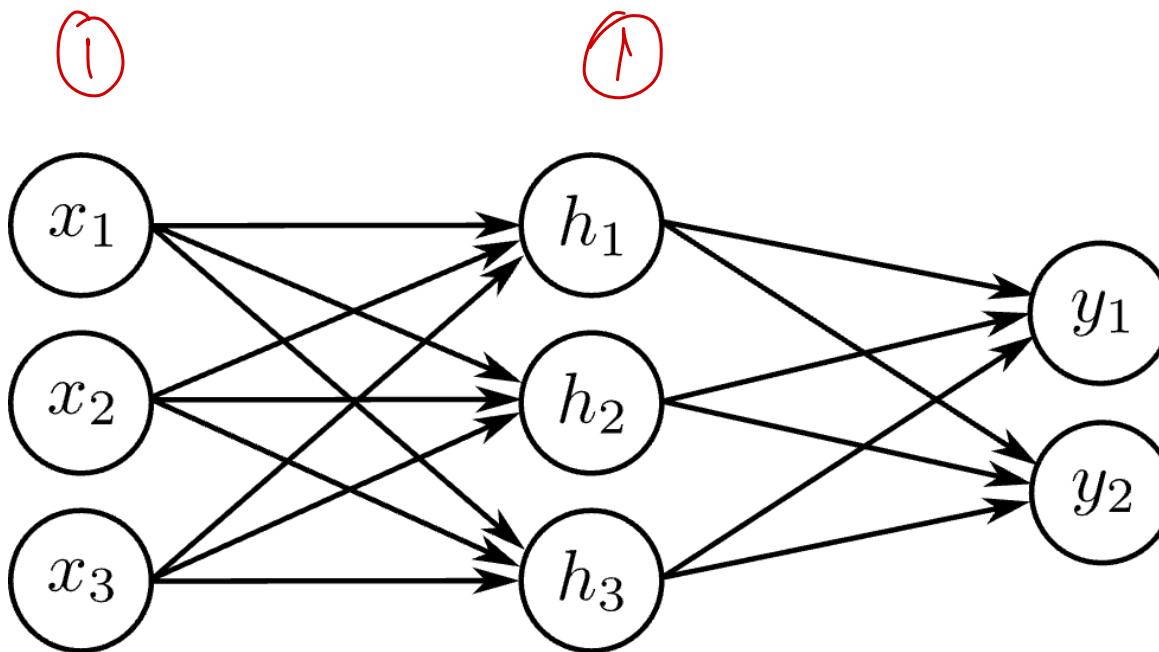
$$h_d = a \left[\theta_{d0} + \sum_{i=1}^{D_i} \theta_{di} x_i \right] \quad y_j = \phi_{j0} + \sum_{d=1}^D \phi_{jd} h_d$$

- e.g., Three inputs, three hidden units, two outputs

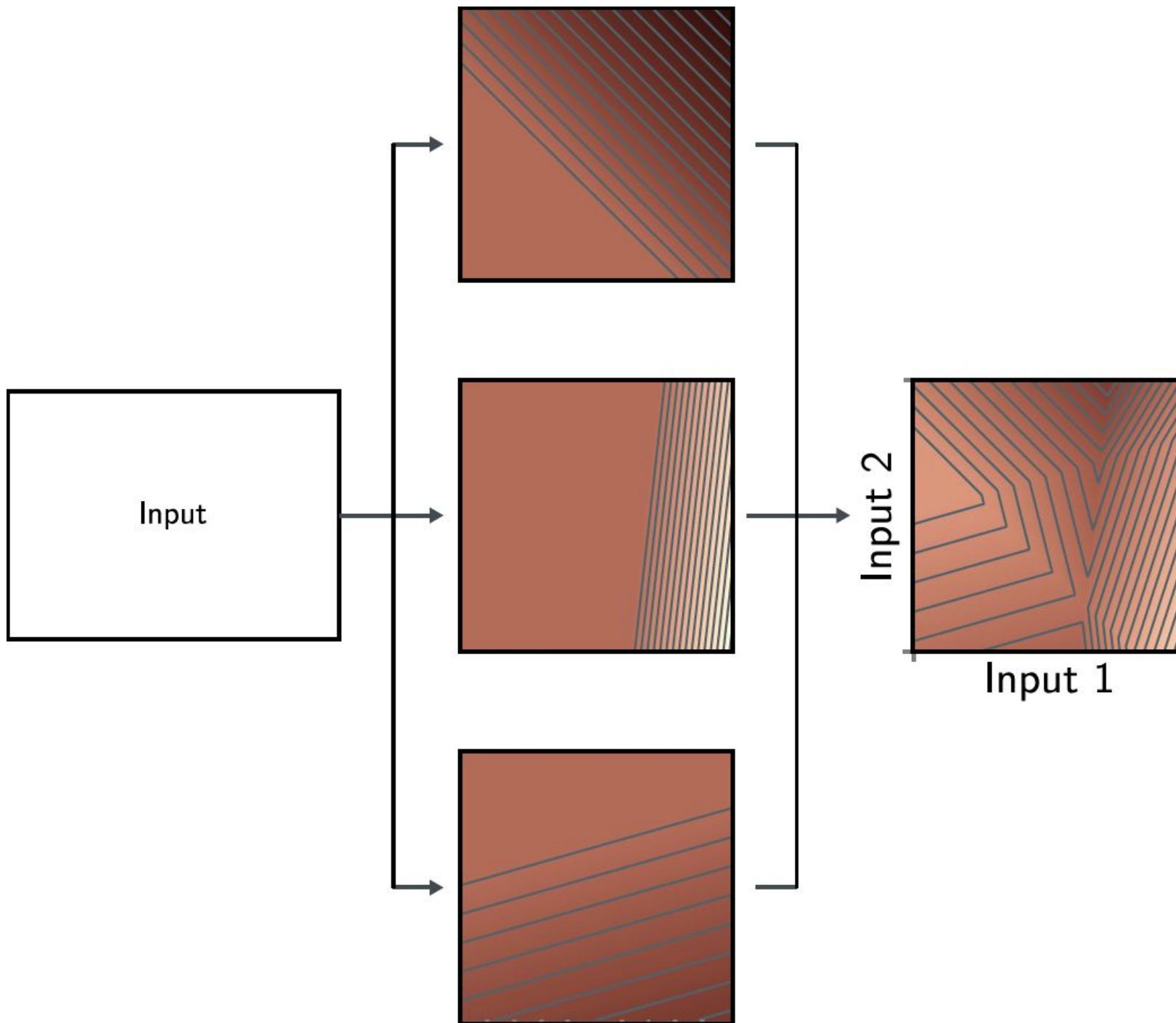


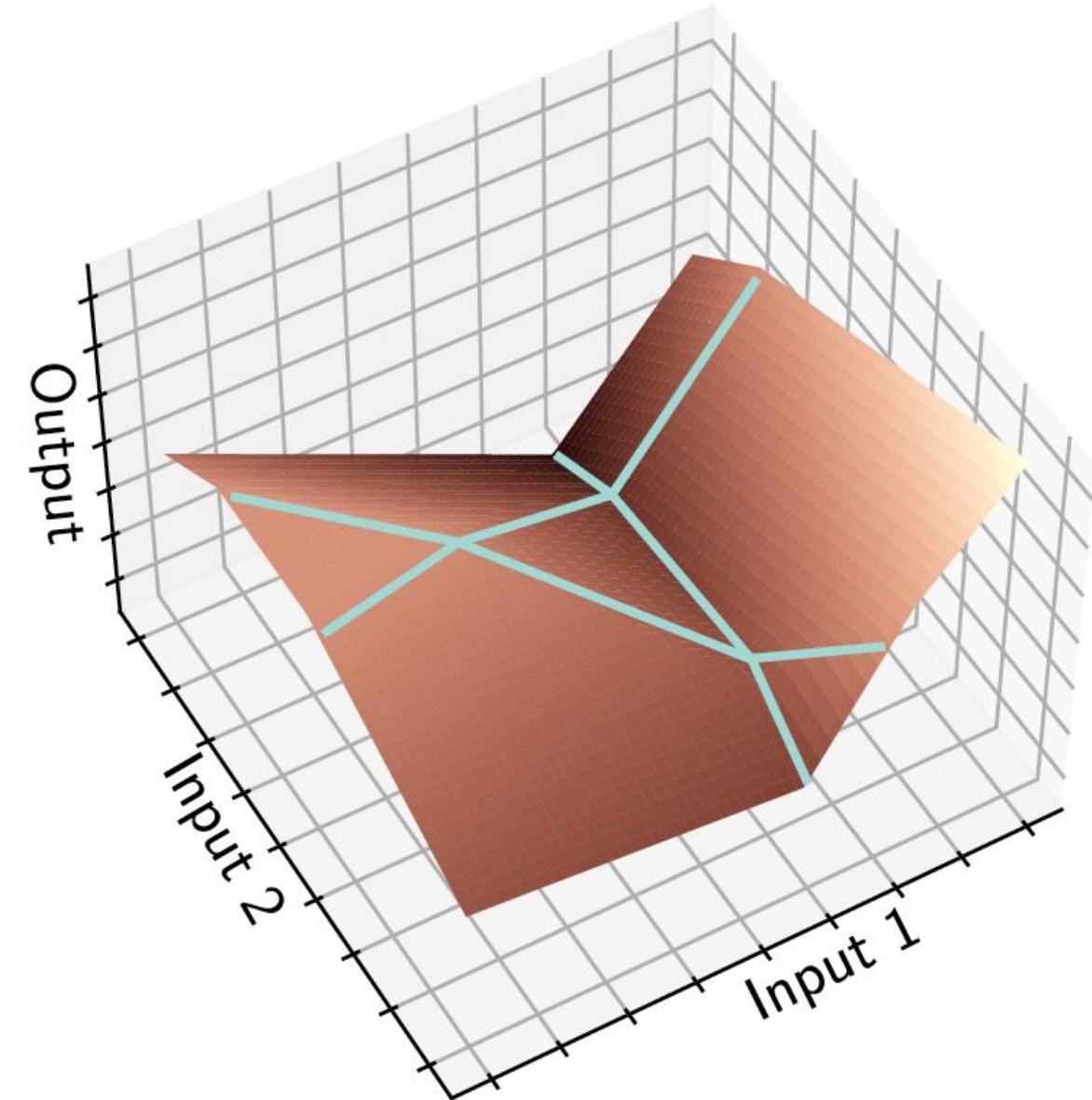
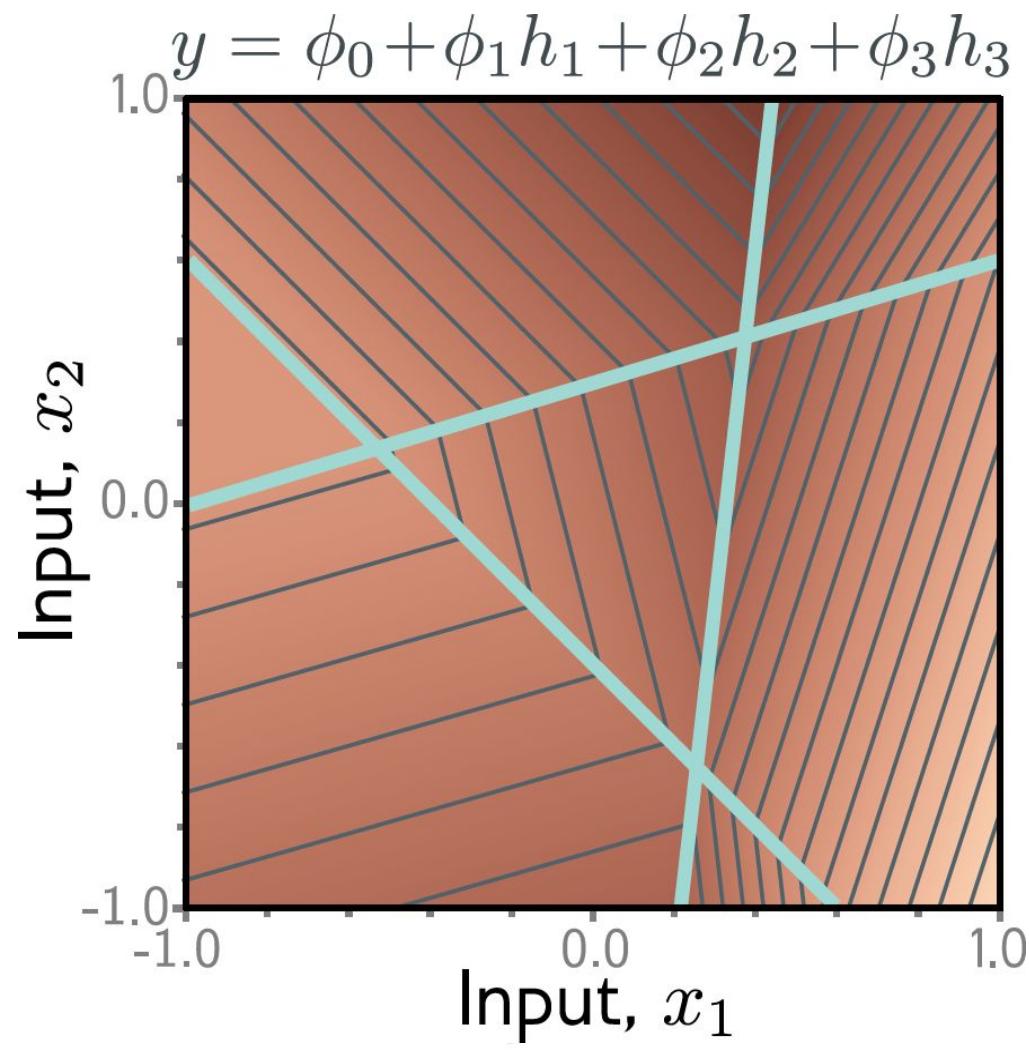
Question:

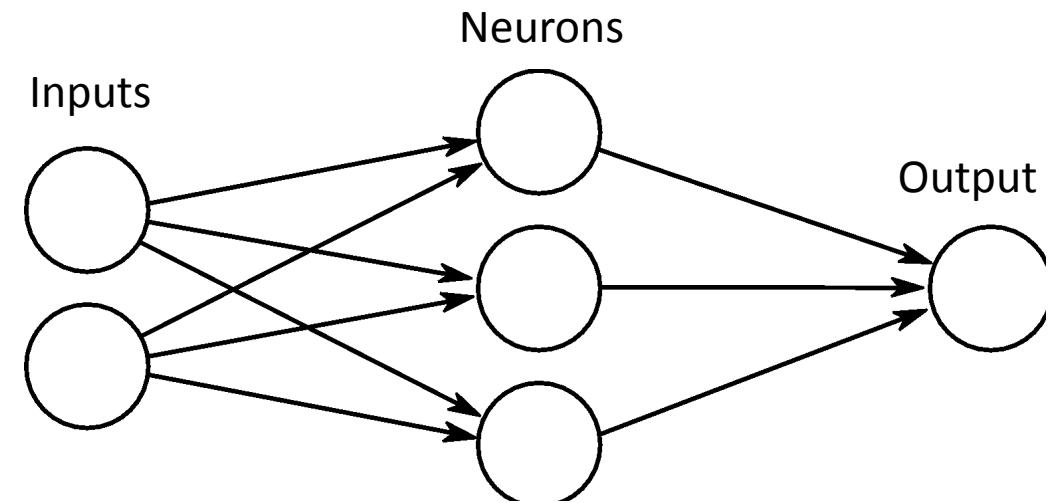
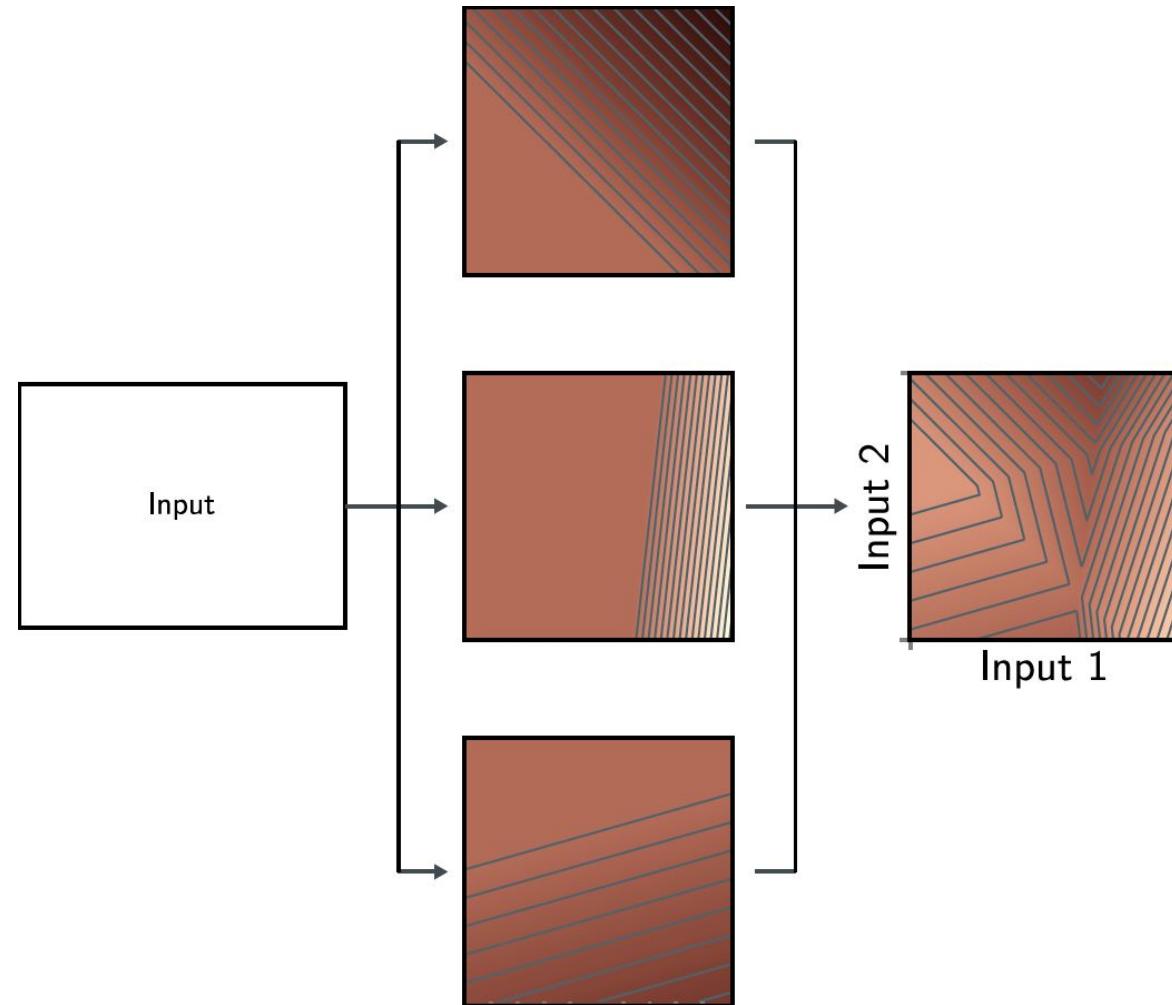
- How many parameters does this model have?



$$\begin{aligned} & 4 \times 3 + 4 \times 2 \\ & = 20 \text{ params} \end{aligned}$$







“neural network”

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

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

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

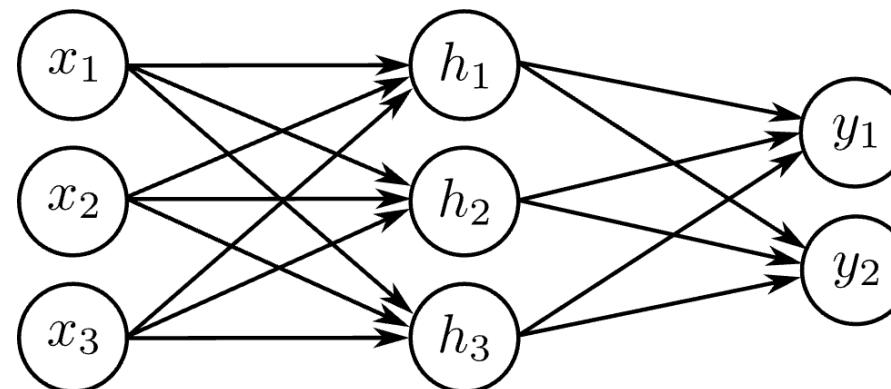
$$h_3 = a[\theta_{30} + \theta_{31}x_1 + \theta_{32}x_2]$$

Arbitrary inputs, hidden units, outputs

- D_o Outputs, D hidden units, and D_i inputs

$$h_d = a \left[\theta_{d0} + \sum_{i=1}^{D_i} \theta_{di} x_i \right] \quad y_j = \phi_{j0} + \sum_{d=1}^D \phi_{jd} h_d$$

- e.g., Three inputs, three hidden units, two outputs

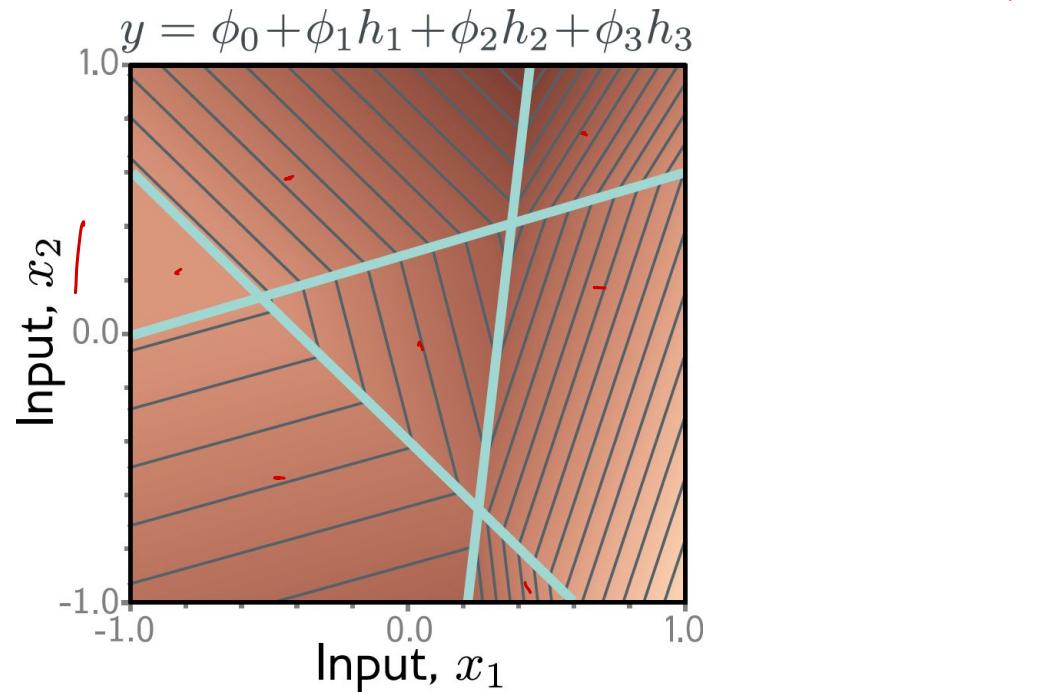


Shallow neural networks

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

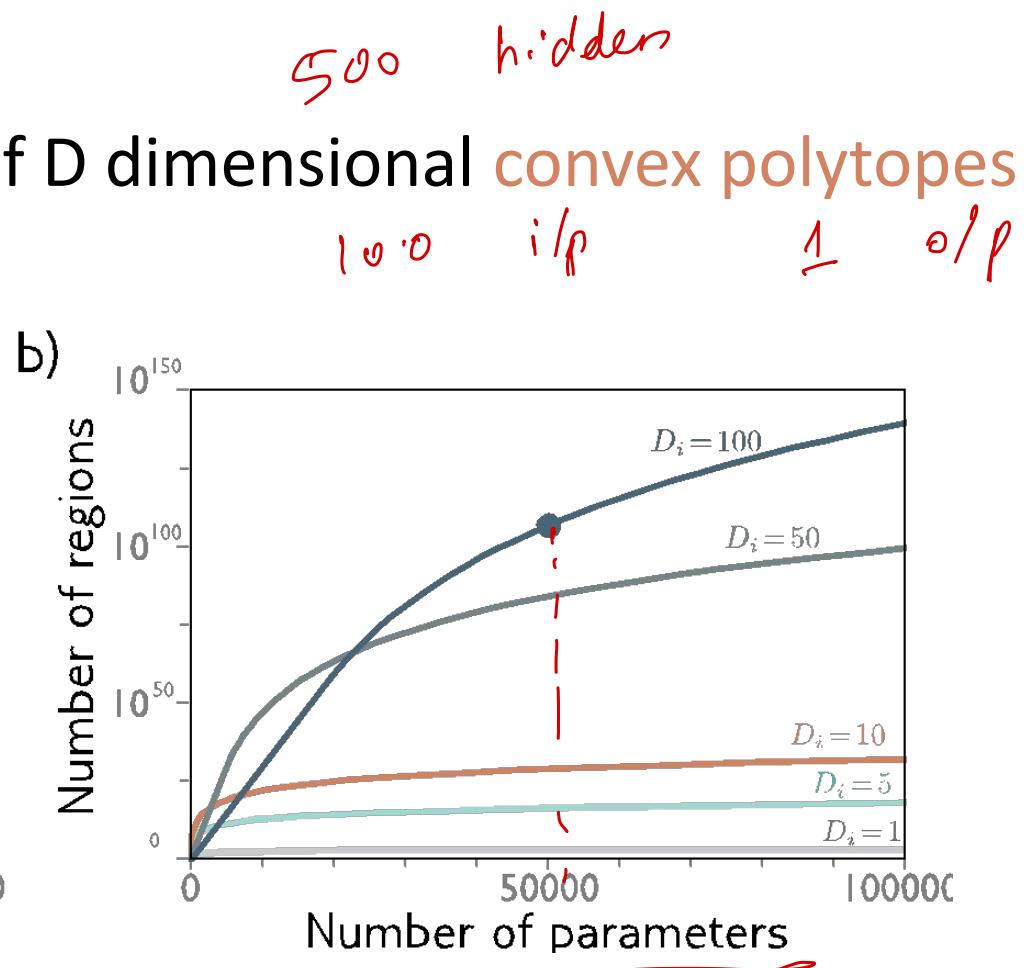
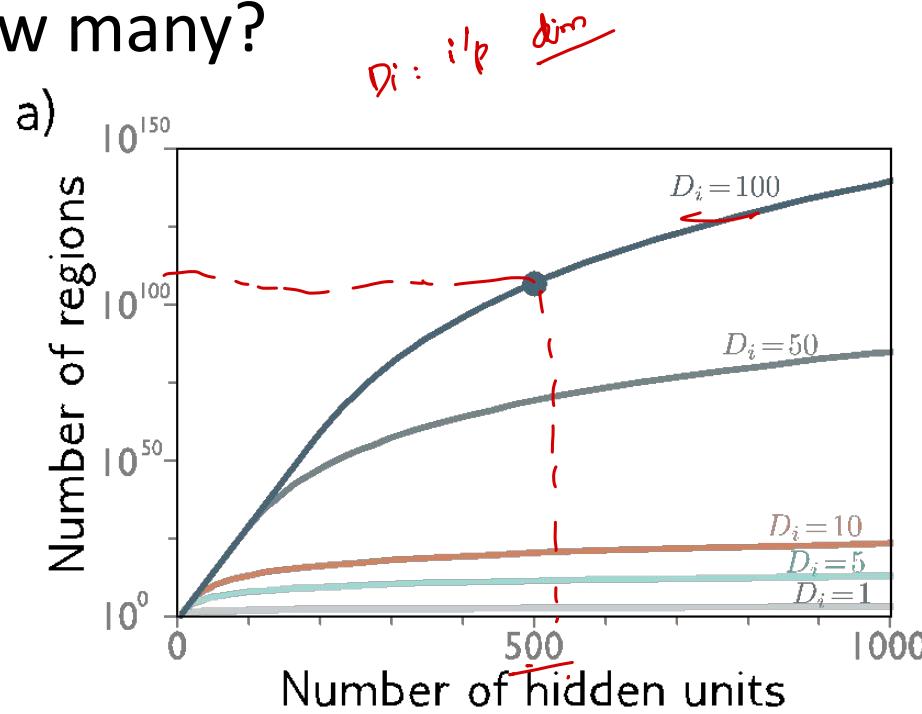
Number of output regions

- In general, each output consists of D dimensional convex polytopes
- With two inputs, and three ~~outputs~~, we saw there were seven polygons:



Number of output regions

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



Highlighted point = 500 hidden units or 51,001 parameters