# Advanced Machine Learning

## M. Ishtiaq

Deep Neural Networks

# Deep neural networks

- Networks with more than one hidden layer

- Intuition becomes more difficult

# Deep neural networks

- Composing two networks
- Combining the two networks into one
- Hyperparameters
- Notation change and general case
- Shallow vs. deep networks

# Composing two networks.

Network 1:

$$h_1 = \mathrm{a}[\theta_{10} + \theta_{11}x]$$
$$h_2 = \mathrm{a}[\theta_{20} + \theta_{21}x] \qquad y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$
$$h_3 = \mathrm{a}[\theta_{30} + \theta_{31}x]$$

Network 2:

$$h_1' = \mathrm{a}[\theta_{10}' + \theta_{11}'y]$$
$$h_2' = \mathrm{a}[\theta_{20}' + \theta_{21}'y] \qquad y' = \phi_0' + \phi_1' h_1' + \phi_2' h_2' + \phi_3' h_3'$$
$$h_3' = \mathrm{a}[\theta_{30}' + \theta_{31}'y]$$

# Composing two networks.

Network 1:

$$h_1 = \mathrm{a}[\theta_{10} + \theta_{11}x]$$
$$h_2 = \mathrm{a}[\theta_{20} + \theta_{21}x]$$
$$h_3 = \mathrm{a}[\theta_{30} + \theta_{31}x]$$
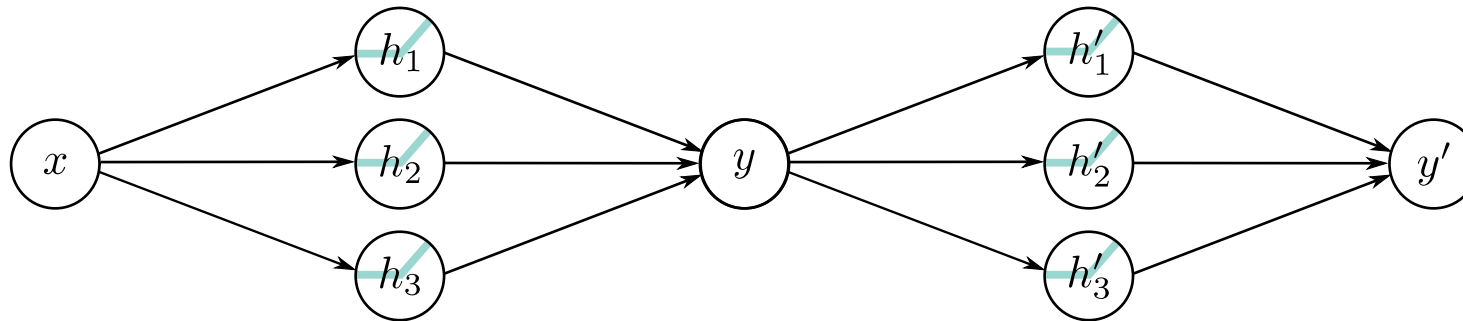
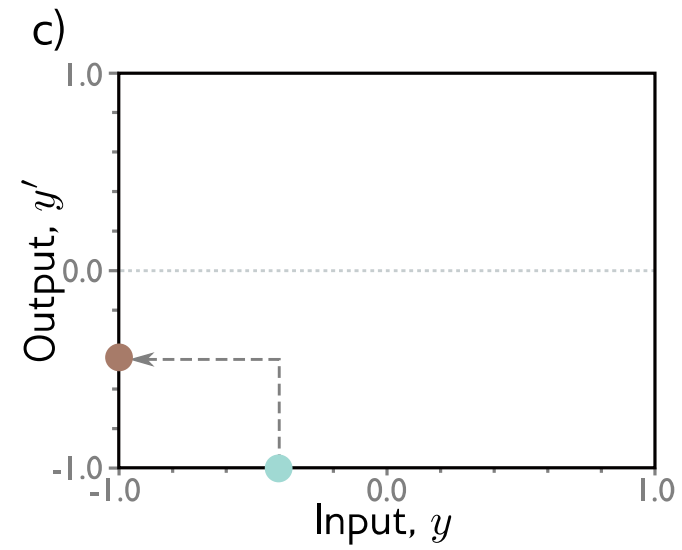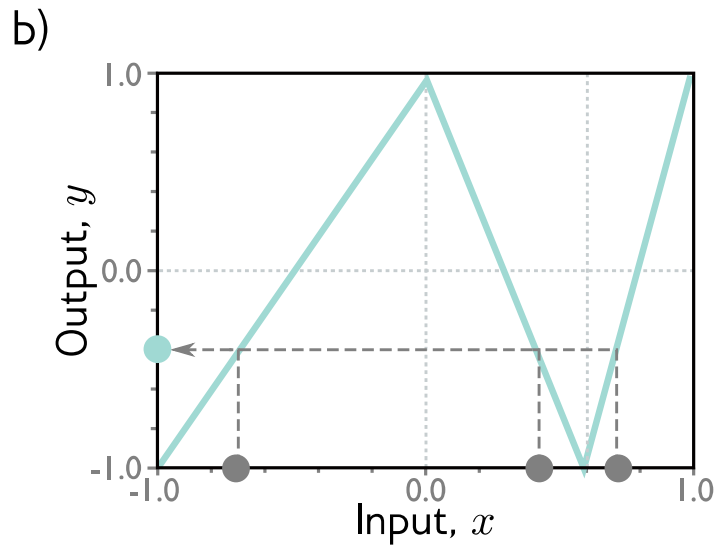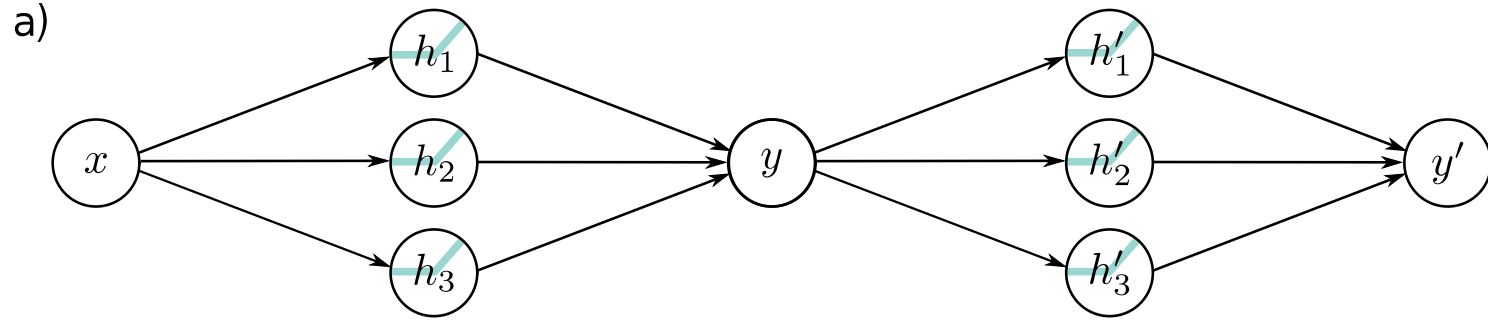$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

Network 2:

$$h'_1 = \mathrm{a}[\theta'_{10} + \theta'_{11}y]$$
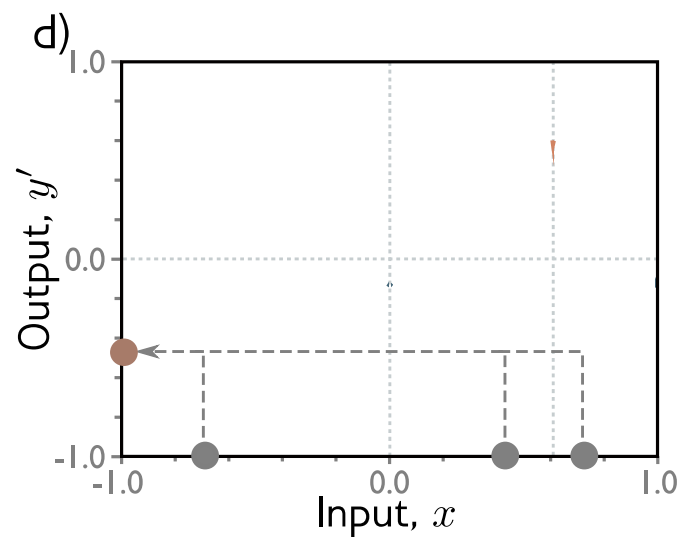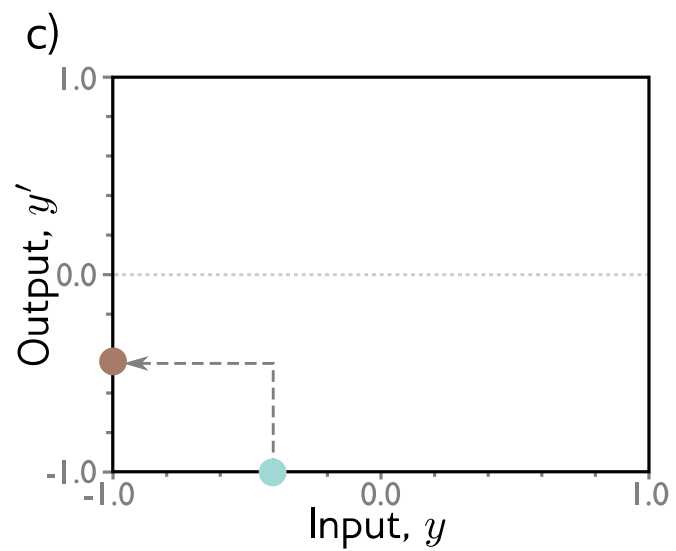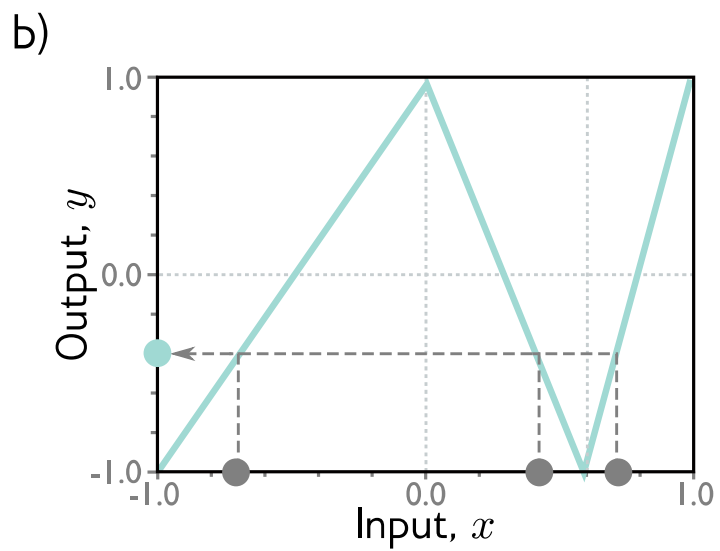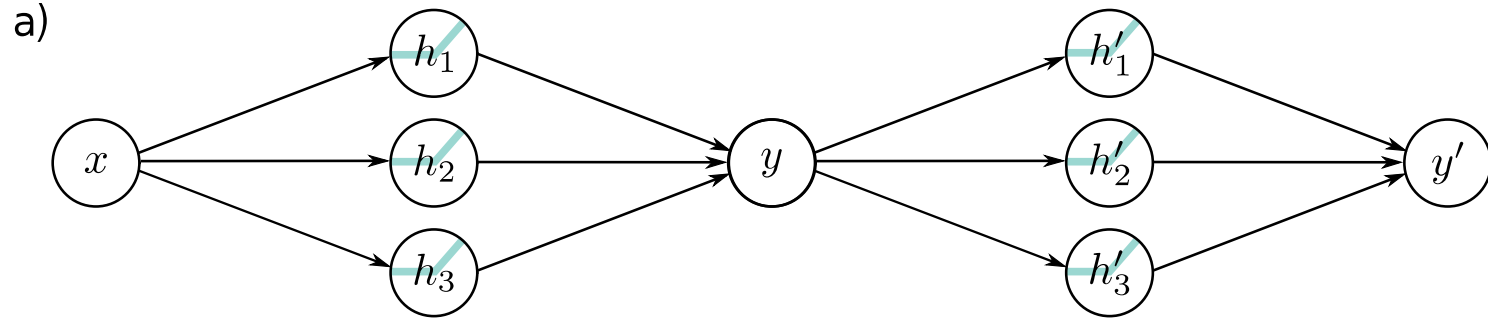$$h'_2 = \mathrm{a}[\theta'_{20} + \theta'_{21}y]$$
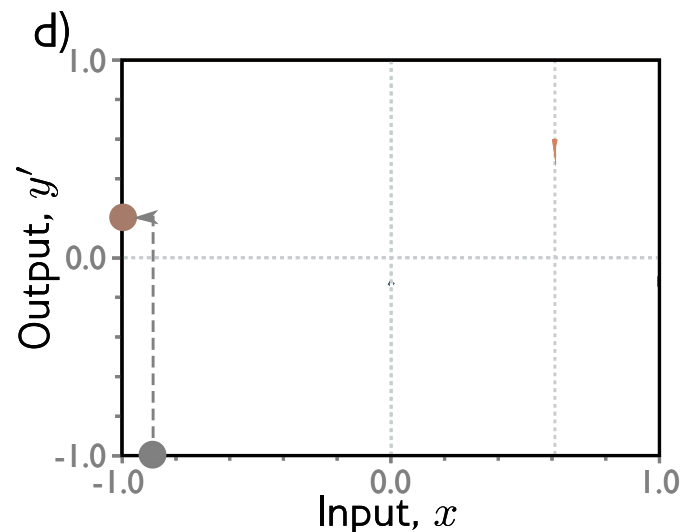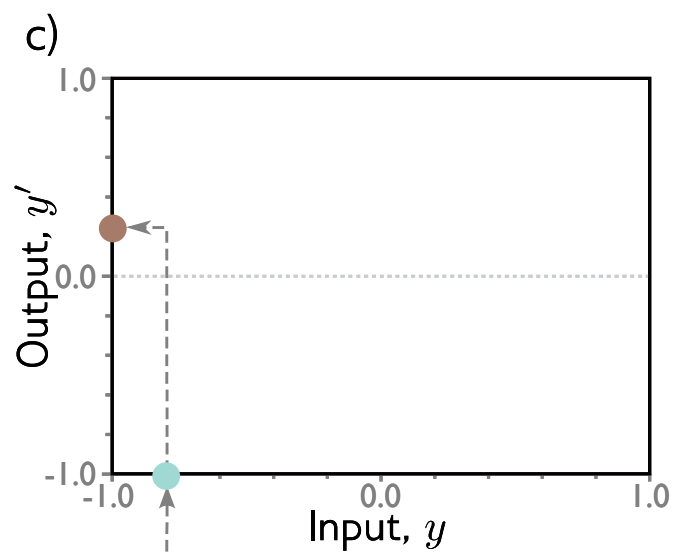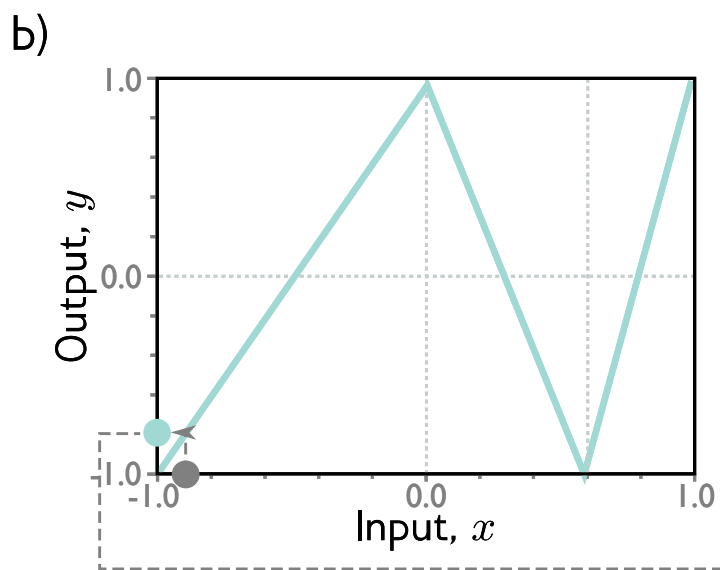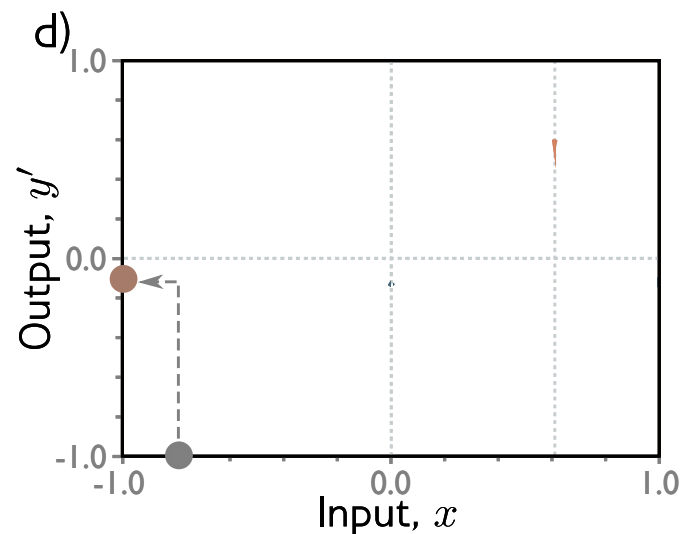$$h'_3 = \mathrm{a}[\theta'_{30} + \theta'_{31}y]$$

$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$

a)

b)

Output, $y$

Input, $x$

c)

Output, $y'$

Input, $y$

a)



b)



c)



d)

a)

b)

c)

d)

13

a)

b)

Output, $y$ vs Input, $x$

c)

Output, $y'$ vs Input, $y$

d)

Output, $y'$ vs Input, $x$

a)



b)



c)



d)

a)



b)



c)



d)

a)



b)



c)



d)

# "Folding analogy"



a) Output, $y$ — Input, $x$

b) Output, $y'$ — Input, $y$

c) Output, $y'$ — Input, $x$

# Comparing to shallow with six hidden units

- 20 parameters
- (at least) 9 regions

- 19 parameters
- Max 7 regions

# Composing networks in 2D



a)

b) Output, $y$

c)

d) Output, $y'$

# Deep neural networks

- Composing two networks
- Combining the two networks into one
- Hyperparameters
- Notation change and general case
- Shallow vs. deep networks

# Combine two networks into one

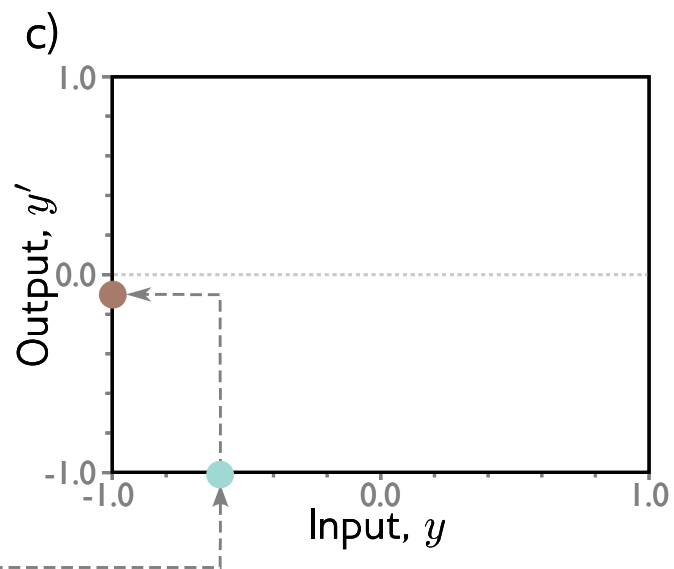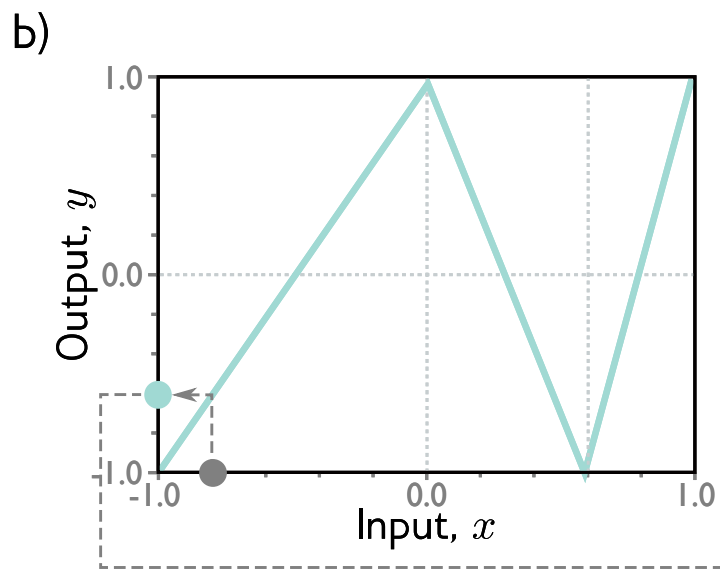$$h_1 = a[\theta_{10} + \theta_{11}x]$$

Network 1:
$$h_2 = a[\theta_{20} + \theta_{21}x] \qquad y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

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

$$h'_1 = a[\theta'_{10} + \theta'_{11}y]$$
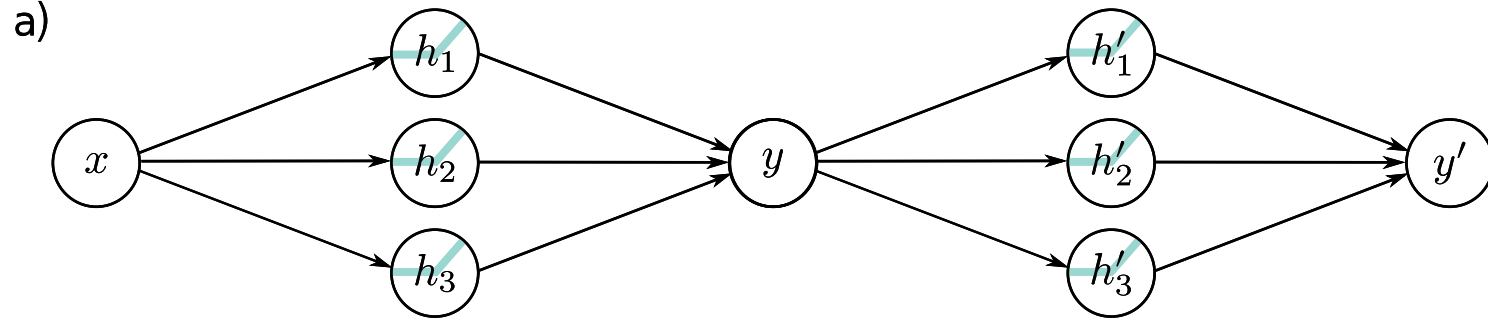
Network 2:
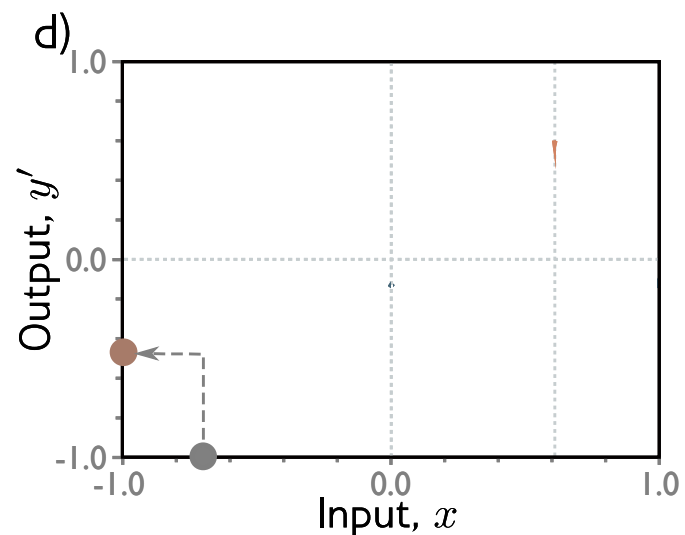$$h'_2 = a[\theta'_{20} + \theta'_{21}y] \qquad y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$

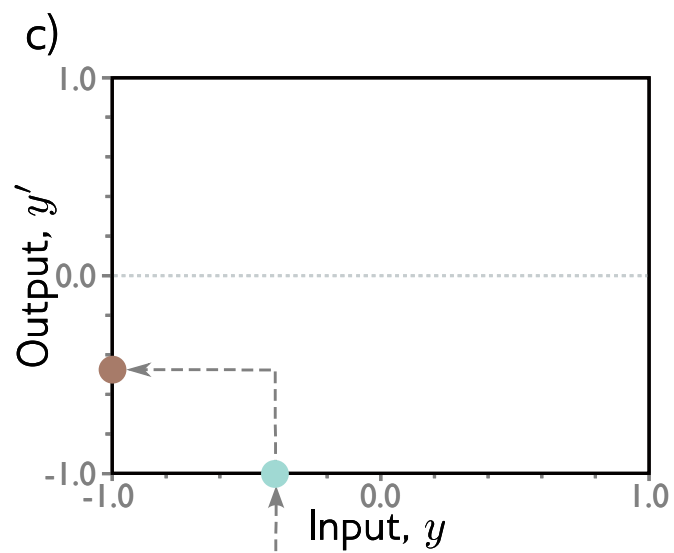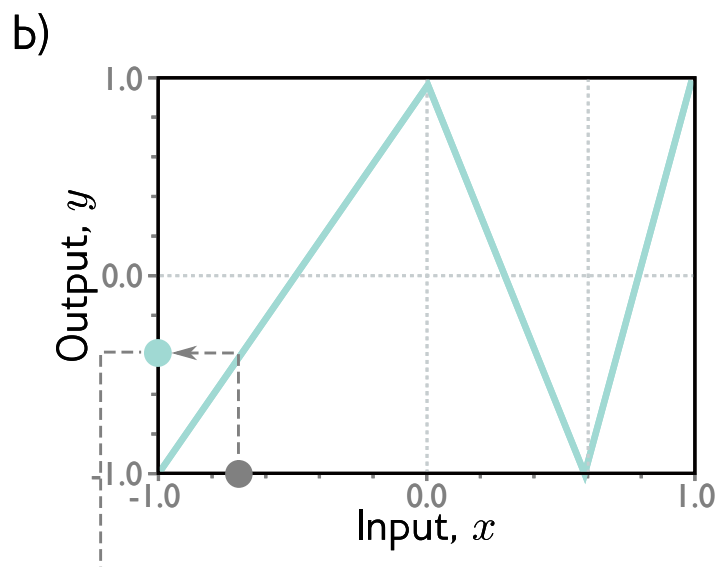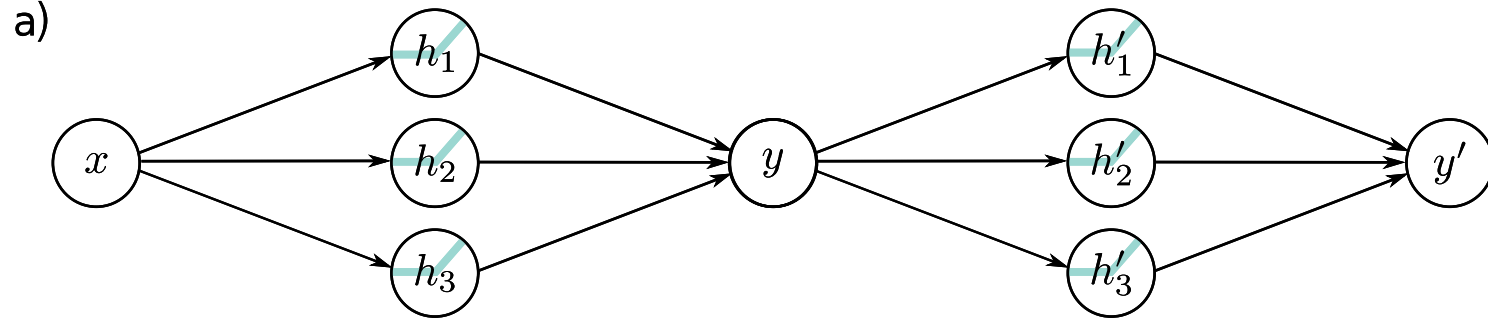$$h'_3 = a[\theta'_{30} + \theta'_{31}y]$$

Hidden units of second network in terms of first:

$$h'_1 = \quad a[\theta'_{10} + \theta'_{11}y] \quad = \quad a[\theta'_{10} + \theta'_{11}\phi_0 + \theta'_{11}\phi_1 h_1 + \theta'_{11}\phi_2 h_2 + \theta'_{11}\phi_3 h_3]$$

$$h'_2 = \quad a[\theta'_{20} + \theta'_{21}y] \quad = \quad a[\theta'_{20} + \theta'_{21}\phi_0 + \theta'_{21}\phi_1 h_1 + \theta'_{21}\phi_2 h_2 + \theta'_{21}\phi_3 h_3]$$

$$h'_3 = \quad a[\theta'_{30} + \theta'_{31}y] \quad = \quad a[\theta'_{30} + \theta'_{31}\phi_0 + \theta'_{31}\phi_1 h_1 + \theta'_{31}\phi_2 h_2 + \theta'_{31}\phi_3 h_3]$$

# Create new variables

$$h_1' = \quad \text{a}[\theta_{10}' + \theta_{11}'y] \quad = \quad \text{a}[\theta_{10}' + \theta_{11}'\phi_0 + \theta_{11}'\phi_1 h_1 + \theta_{11}'\phi_2 h_2 + \theta_{11}'\phi_3 h_3]$$

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

$$h_3' = \quad \text{a}[\theta_{30}' + \theta_{31}'y] \quad = \quad \text{a}[\theta_{30}' + \theta_{31}'\phi_0 + \theta_{31}'\phi_1 h_1 + \theta_{31}'\phi_2 h_2 + \theta_{31}'\phi_3 h_3]$$

$$h_1' = \text{a}[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

$$h_2' = \text{a}[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$

$$h_3' = \text{a}[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

# Two-layer network

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

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

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

$$h_1' = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

$$h_2' = a[\psi_{20} + \psi_{21}h_2 + \psi_{22}h_2 + \psi_{23}h_3]$$

$$h_3' = a[\psi_{30} + \psi_{31}h_2 + \psi_{32}h_2 + \psi_{33}h_3]$$

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

# Two-layer network as one equation

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

$$h_1' = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$
$$h_2' = a[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$
$$h_3' = a[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

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

$$y' = \phi_0' + \phi_1'a\left[\psi_{10} + \psi_{11}a[\theta_{10} + \theta_{11}x] + \psi_{12}a[\theta_{20} + \theta_{21}x] + \psi_{13}a[\theta_{30} + \theta_{31}x]\right]$$
$$+ \phi_2'a[\psi_{20} + \psi_{21}a[\theta_{10} + \theta_{11}x] + \psi_{22}a[\theta_{20} + \theta_{21}x] + \psi_{23}a[\theta_{30} + \theta_{31}x]]$$
$$+ \phi_3'a[\psi_{30} + \psi_{31}a[\theta_{10} + \theta_{11}x] + \psi_{32}a[\theta_{20} + \theta_{21}x] + \psi_{33}a[\theta_{30} + \theta_{31}x]]$$

# Remember shallow network with two outputs?

- 1 input, 4 hidden units, 2 outputs

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

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

# Networks as composing functions

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

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

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

$$h_1' = \mathrm{a}[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

$$h_2' = \mathrm{a}[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$

$$h_3' = \mathrm{a}[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

Consider the pre-activations at the second hidden units
At this point, it's a one--layer network with three outputs

# Networks as composing functions

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

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

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

$$h'_1 = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

$$h'_2 = a[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$

$$h'_3 = a[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

Consider the pre-activations at the second hidden units
At this point, it's a one--layer network with three outputs

a) $\psi_{10}+\psi_{11}h_1+\psi_{12}h_2+\psi_{13}h_3$

b) $\psi_{20}+\psi_{21}h_1+\psi_{22}h_2+\psi_{23}h_3$

c) $\psi_{30}+\psi_{31}h_1+\psi_{32}h_2+\psi_{33}h_3$

a) $\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3$

b) $\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3$

c) $\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3$

d) $h'_1 = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$

e) $h'_2 = a[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$

f) $h'_3 = a[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$

d) $h'_1 =$
$\mathsf{a}[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$

e) $h'_2 =$
$\mathsf{a}[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$

f) $h'_3 =$
$\mathsf{a}[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$

g) $\phi'_1 h'_1$

h) $\phi'_2 h'_2$

i) $\phi'_3 h'_3$

g) $\phi_1' h_1'$

h) $\phi_2' h_2'$

i) $\phi_3' h_3'$

j) $\phi_0' + \phi_1' h_1' + \phi_2' h_2' + \phi_3' h_3'$

# Deep neural networks

- Composing two networks
- Combining the two networks into one
- Hyperparameters
- Notation change and general case
- Shallow vs. deep networks

# Hyperparameters

- K layers = depth of network
- $D_k$ hidden units per layer = width of network

- These are called hyperparameters – chosen before training the network
- Can try retraining with different hyperparameters – hyperparameter optimization or hyperparameter search

# Deep neural networks

- Composing two networks
- Combining the two networks into one
- Hyperparameters
- Notation change and general case
- Shallow vs. deep networks

# Notation change #1

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

$$h'_1 = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$
$$h'_2 = a[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$
$$h'_3 = a[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$

# Notation change #1

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

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = \mathbf{a} \left[ \begin{bmatrix} \theta_{10} \\ \theta_{20} \\ \theta_{30} \end{bmatrix} + \begin{bmatrix} \theta_{11} \\ \theta_{21} \\ \theta_{31} \end{bmatrix} x \right]$$
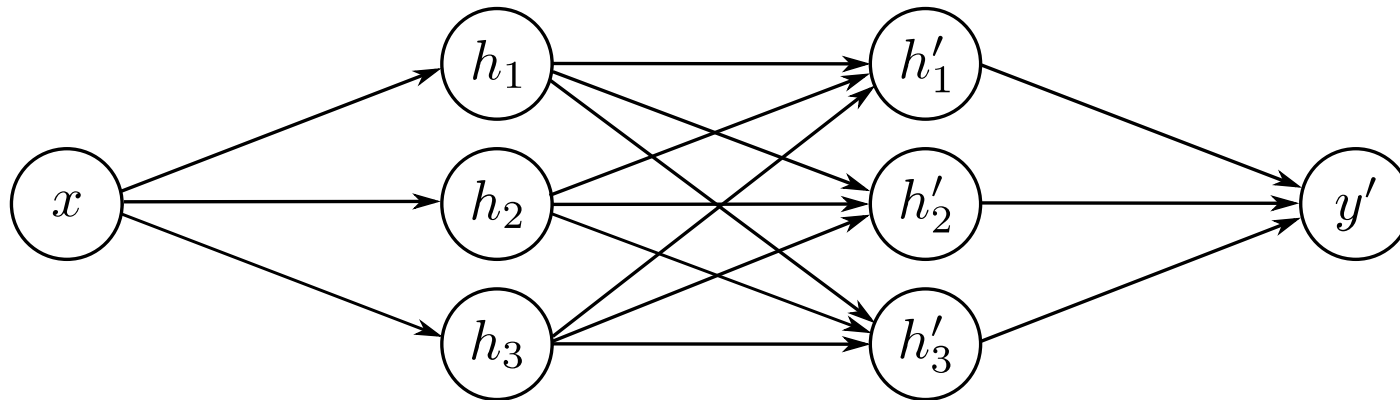
$$h_1' = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$
$$h_2' = a[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$
$$h_3' = a[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

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

# Notation change #1

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

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

$$h'_1 = \mathrm{a}[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$
$$h'_2 = \mathrm{a}[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$
$$h'_3 = \mathrm{a}[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

$$\begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix} = \mathbf{a} \left[ \begin{bmatrix} \psi_{10} \\ \psi_{20} \\ \psi_{30} \end{bmatrix} + \begin{bmatrix} \psi_{11} & \psi_{12} & \psi_{13} \\ \psi_{21} & \psi_{22} & \psi_{23} \\ \psi_{31} & \psi_{32} & \psi_{33} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} \right]$$

# Notation change #1

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

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = a \left[ \begin{bmatrix} \theta_{10} \\ \theta_{20} \\ \theta_{30} \end{bmatrix} + \begin{bmatrix} \theta_{11} \\ \theta_{21} \\ \theta_{31} \end{bmatrix} x \right]$$
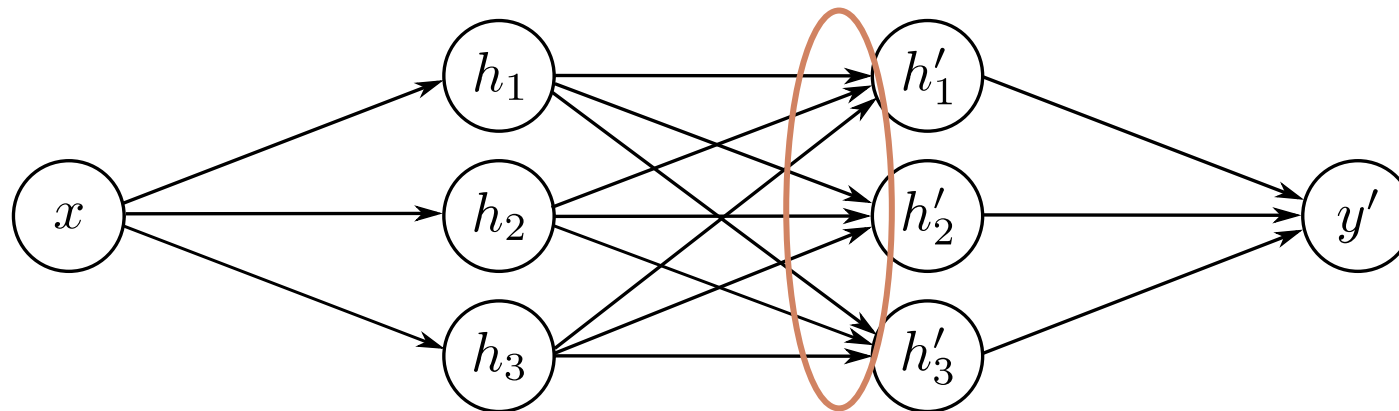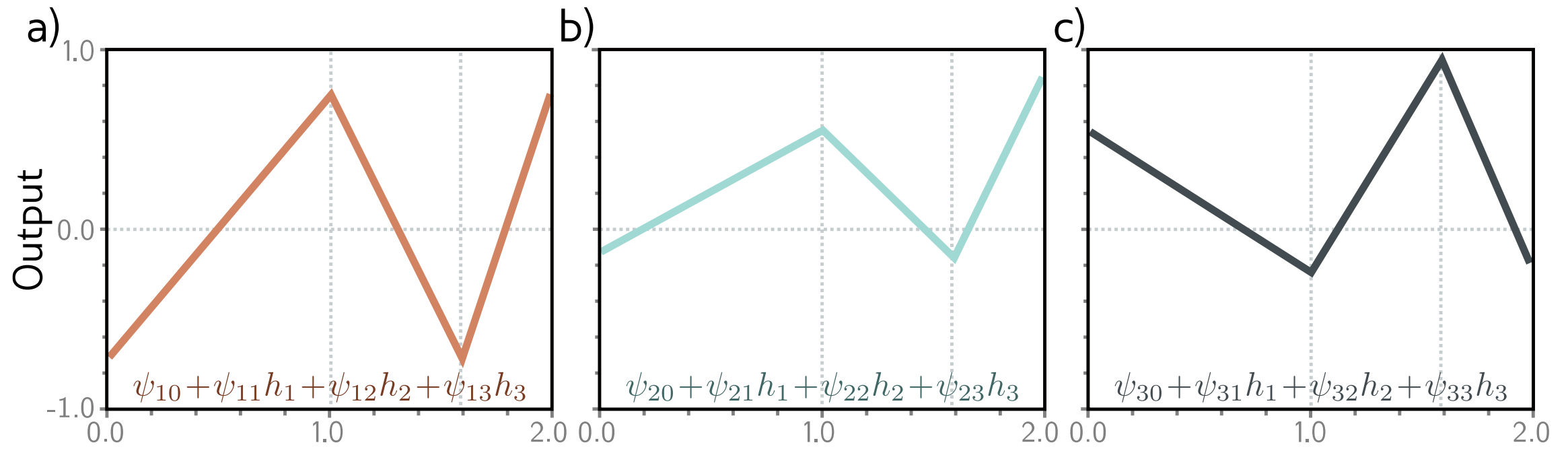
$$h'_1 = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$
$$h'_2 = a[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$
$$h'_3 = a[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

$$\begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix} = a \left[ \begin{bmatrix} \psi_{10} \\ \psi_{20} \\ \psi_{30} \end{bmatrix} + \begin{bmatrix} \psi_{11} & \psi_{12} & \psi_{13} \\ \psi_{21} & \psi_{22} & \psi_{23} \\ \psi_{32} & \psi_{32} & \psi_{33} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} \right]$$

$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$

$$y' = \phi'_0 + \begin{bmatrix} \phi'_1 & \phi'_2 & \phi'_3 \end{bmatrix} \begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix}$$

# Notation change #2

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = \mathbf{a} \left[ \begin{bmatrix} \theta_{10} \\ \theta_{20} \\ \theta_{30} \end{bmatrix} + \begin{bmatrix} \theta_{11} \\ \theta_{21} \\ \theta_{31} \end{bmatrix} x \right] \longrightarrow \mathbf{h} = \mathbf{a} \left[ \boldsymbol{\theta}_0 + \boldsymbol{\theta} x \right]$$

$$\begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix} = \mathbf{a} \left[ \begin{bmatrix} \psi_{10} \\ \psi_{20} \\ \psi_{30} \end{bmatrix} + \begin{bmatrix} \psi_{11} & \psi_{12} & \psi_{13} \\ \psi_{21} & \psi_{22} & \psi_{23} \\ \psi_{32} & \psi_{32} & \psi_{33} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} \right] \longrightarrow \mathbf{h}' = \mathbf{a} \left[ \boldsymbol{\psi}_0 + \boldsymbol{\Psi} \mathbf{h} \right]$$

$$y' = \phi'_0 + \begin{bmatrix} \phi'_1 & \phi'_2 & \phi'_3 \end{bmatrix} \begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix} \longrightarrow y = \phi'_0 + \boldsymbol{\phi}' \mathbf{h}'$$

# Notation change #3

$$\mathbf{h} = \mathbf{a}\left[\boldsymbol{\theta}_0 + \boldsymbol{\theta}x\right] \longrightarrow \mathbf{h}_1 = \mathbf{a}[\boldsymbol{\beta}_0 + \boldsymbol{\Omega}_0\mathbf{x}]$$

$$\mathbf{h}' = \mathbf{a}\left[\boldsymbol{\psi}_0 + \boldsymbol{\Psi}\mathbf{h}\right] \longrightarrow \mathbf{h}_2 = \mathbf{a}[\boldsymbol{\beta}_1 + \boldsymbol{\Omega}_1\mathbf{h}_1]$$

$$y = \boldsymbol{\phi}'_0 + \boldsymbol{\phi}'\mathbf{h}' \longrightarrow \mathbf{y} = \boldsymbol{\beta}_2 + \boldsymbol{\Omega}_2\mathbf{h}_2$$

# Notation change #3

$$\mathbf{h} = \mathbf{a}\left[\boldsymbol{\theta}_0 + \boldsymbol{\theta}x\right]$$

$$\mathbf{h}_1 = \mathbf{a}[\boldsymbol{\beta}_0 + \boldsymbol{\Omega}_0\mathbf{x}]$$

$$\mathbf{h}' = \mathbf{a}\left[\boldsymbol{\psi}_0 + \boldsymbol{\Psi}\mathbf{h}\right]$$

$$\mathbf{h}_2 = \mathbf{a}[\boldsymbol{\beta}_1 + \boldsymbol{\Omega}_1\mathbf{h}_1]$$

$$y = \phi_0' + \phi'\mathbf{h}'$$

$$\mathbf{y} = \boldsymbol{\beta}_2 + \boldsymbol{\Omega}_2\mathbf{h}_2$$

# General equations for deep network

$$\mathbf{h}_1 = \mathbf{a}[\boldsymbol{\beta}_0 + \boldsymbol{\Omega}_0 \mathbf{x}]$$

$$\mathbf{h}_2 = \mathbf{a}[\boldsymbol{\beta}_1 + \boldsymbol{\Omega}_1 \mathbf{h}_1]$$

$$\mathbf{h}_3 = \mathbf{a}[\boldsymbol{\beta}_2 + \boldsymbol{\Omega}_2 \mathbf{h}_2]$$

$$\vdots$$

$$\mathbf{h}_K = \mathbf{a}[\boldsymbol{\beta}_{K-1} + \boldsymbol{\Omega}_{K-1} \mathbf{h}_{K-1}]$$

$$\mathbf{y} = \boldsymbol{\beta}_K + \boldsymbol{\Omega}_K \mathbf{h}_K,$$

---

$$\mathbf{y} = \boldsymbol{\beta}_K + \boldsymbol{\Omega}_K \mathbf{a}\left[\boldsymbol{\beta}_{K-1} + \boldsymbol{\Omega}_{K-1}\mathbf{a}\left[\ldots \boldsymbol{\beta}_2 + \boldsymbol{\Omega}_2\mathbf{a}\left[\boldsymbol{\beta}_1 + \boldsymbol{\Omega}_1\mathbf{a}\left[\boldsymbol{\beta}_0 + \boldsymbol{\Omega}_0\mathbf{x}\right]\right]\ldots\right]\right]$$

# Example



$\boldsymbol{\beta}_0 \in \mathbb{R}^4$   $\boldsymbol{\beta}_1 \in \mathbb{R}^2$   $\boldsymbol{\beta}_2 \in \mathbb{R}^3$   $\boldsymbol{\beta}_3 \in \mathbb{R}^2$

$\boldsymbol{\Omega}_0 \in \mathbb{R}^{4 \times 3}$   $\boldsymbol{\Omega}_1 \in \mathbb{R}^{2 \times 4}$   $\boldsymbol{\Omega}_2 \in \mathbb{R}^{3 \times 2}$   $\boldsymbol{\Omega}_3 \in \mathbb{R}^{2 \times 3}$

Input, $\mathbf{x}$    Hidden layer, $\mathbf{h}_1$    Hidden layer, $\mathbf{h}_2$    Hidden layer, $\mathbf{h}_3$    Output, $\mathbf{y}$

$D_i = 3$    $D_1 = 4$    $D_2 = 2$    $D_3 = 3$    $D_o = 2$

# Deep neural networks

- Composing two networks
- Combining the two networks into one
- Hyperparameters
- Notation change and general case
- Shallow vs. deep networks

# Shallow vs. deep networks

The best results are created by deep networks with many layers.

- 50-1000 layers for most applications
- Best results in
    - Computer vision
    - Natural language processing
    - Graph neural networks
    - Generative models
    - Reinforcement learning

All use deep networks.
But why?

# Shallow vs. deep networks

1. Ability to approximate different functions?


Both obey the universal approximation theorem.


Argument: One layer is enough, and for deep networks could arrange for the other layers to compute the identity function.

# Shallow vs. deep networks

2. Number of linear regions per parameter

# Number of linear regions per parameter



a) Input dimension $D_i = 1$

$K=5$
$K=4$
$K=3$
$K=2$
$K=1$

Number of regions

Number of parameters

5 layers
10 hidden units per layer
471 parameters
161,501 linear regions

# Number of linear regions per parameter



a) Input dimension $D_i = 1$

Number of regions / Number of parameters

$K=5$, $K=4$, $K=3$, $K=2$, $K=1$

5 layers
10 hidden units per layer
471 parameters
161,501 linear regions

b) Input dimension $D_i = 10$

Number of regions / Number of parameters

$K=5$, $K=4$, $K=3$, $K=2$, $K=1$

5 layers
50 hidden units per layer
10,801 parameters
$>10^{40}$ linear regions

# Shallow vs. deep networks

2. Number of linear regions per parameter

- Deep networks create many more regions per parameters
- But there are dependencies between them
    - Think of folding example
    - Perhaps similar symmetries in real-world functions? Unknown

# Shallow vs. Deep Networks

3. Depth efficiency

- There are some functions that require a shallow network with exponentially more hidden units than a deep network to achieve an equivalent approximation

- This is known as the depth efficiency of deep networks

- But do the real-world functions we want to approximate have this property?  Unknown.

# Shallow vs. Deep Networks

4. Large structured networks

- Think about images as input – might be 1M pixels
- Fully connected works not practical
- Answer is to have weights that only operate locally, and share across image
- This leads to convolutional networks
- Gradually integrate information from across the image – needs multiple layers
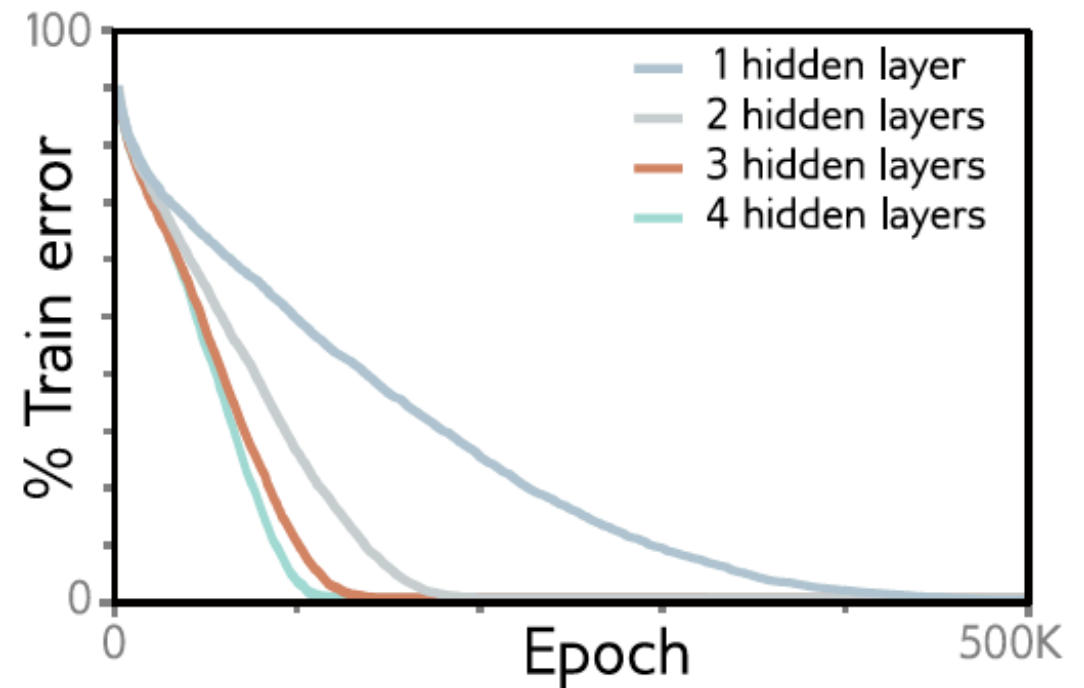
# Shallow vs. Deep Networks

5. Fitting and generalization

- Fitting of deep models seems to be easier up to about 20 layers
- Then needs various tricks to train deeper networks, so (in vanilla form), fitting becomes harder
- Generalization is good in deep networks. Why?

# Shallow vs. Deep Networks

## 5. Fitting and generalization

**Figure 20.2** MNIST-1D training. Four fully connected networks were fit to 4000 MNIST-1D examples with random labels using full batch gradient descent, He initialization, no momentum or regularization, and learning rate 0.0025. Models with 1,2,3,4 layers had 298, 100, 75, and 63 hidden units per layer and 15208, 15210, 15235, and 15139 parameters, respectively. All models train successfully, but deeper models require fewer epochs.

# Where are we going?

- We have defined families of very flexible networks that map multiple inputs to multiple outputs

- Now we need to train them
    - How to choose loss functions
    - How to find minima of the loss function
    - How to do this in particular for deep networks

- Then we need to test them