

# Medical Diagnostics

Tutorial 2: Deep Neural Networks

Spring 2025

Faculty of Biotechnology and Food Engineering  
Technion Israel Institute of Technology

# Table of contents

- 1 Intro
- 2 Tutorial 1 recap
- 3 Perceptron → Multi-Layered Perceptron (MLP)
- 4 Backpropagation
- 5 Implementing and training models

# 1 Intro

# 1.1 Ongoing security situation

התראות חיירום

**HW1 OUT NOW**

00660121 - Medical Diagnostics



# 1.2 Course updates

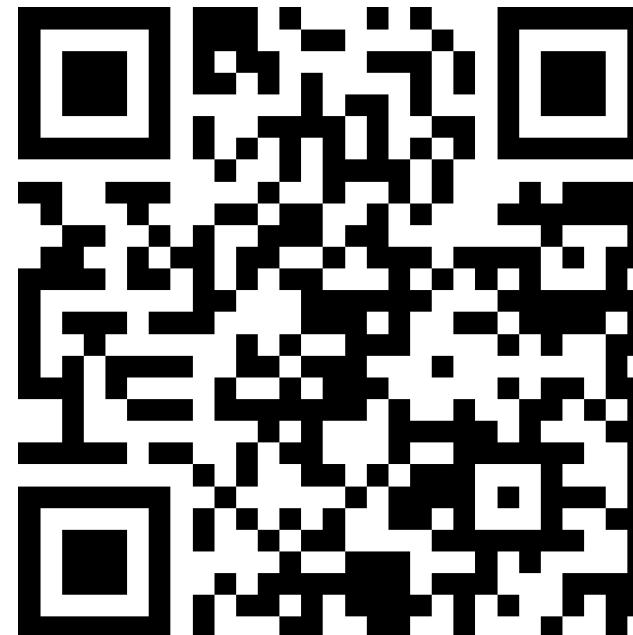
- new HW1 due date is July 10th
  - 10% of final grade
  - Use the **forum**, read existing posts
- HW2+Final project merged
  - 25% of final grade
  - Released on July 10th
  - Due date: September 7th
- Workshops (HW) are on Monday (15:30-17:30) and Wednesday (10:00-11:30)
- Office hours (course material) can be scheduled through mail

# 1.3 Slido

Open [app.sli.do](https://app.sli.do) on your browser

Enter code **00660121**

Answer and ask questions in real-time



00660121 - Medical Diagnostics

# Example: How is Homework 1 going?

# 2 Tutorial 1 recap

# 2.1 Basic notations of ML

Symbol

---

$\mathcal{X}$

---

$\mathcal{Y}$

---

$D = \{(x_i, y_i)\}_{i=1}^N$

---

$f : \mathcal{X} \rightarrow \mathcal{Y}$

---

$\mathcal{A}$

---

$\mathcal{H}$

Symbol

---

$h_{\theta^*} \approx f$

---

$\mathcal{L}(h_{\theta})$

---

$\theta$

---

$\eta$

---

$\nabla_{\theta} \mathcal{L}(h_{\theta})$

# Model training workflow:

1. Prepare data  $D$
2. Choose model  $\mathcal{H}$
3. Train model: Find  $h_{\theta^*}$  s.t.  $\mathcal{L}(h_{\theta^*}) = \min_{h_{\theta} \in \mathcal{H}} \mathcal{L}(h_{\theta})$ 
  - Loss function
  - Optimization algorithm
4. Tune hyperparameters: Use validation set
5. Evaluate: Assess performance on test set
6. Deploy: Use trained model ( $h_{\theta}$ ) for predictions

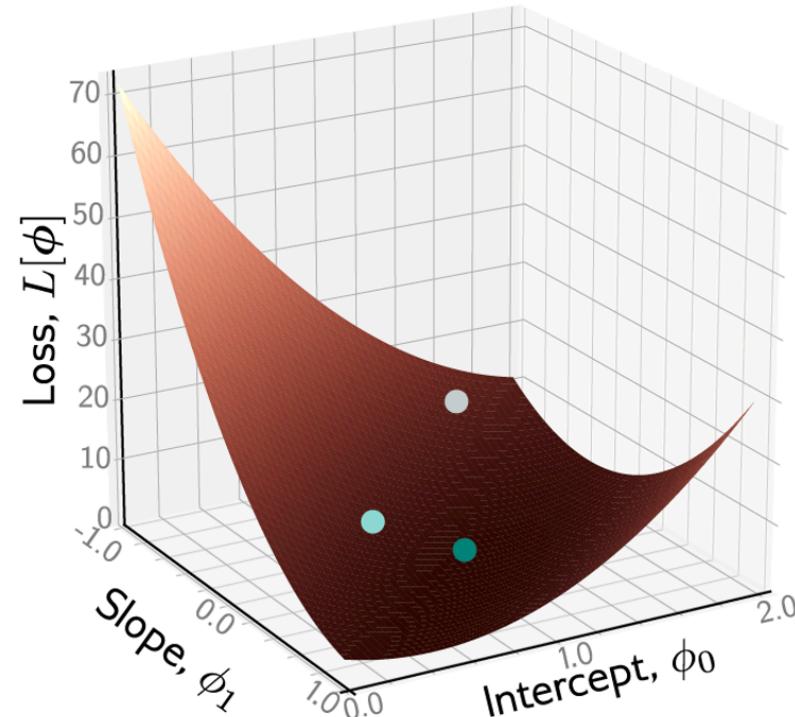
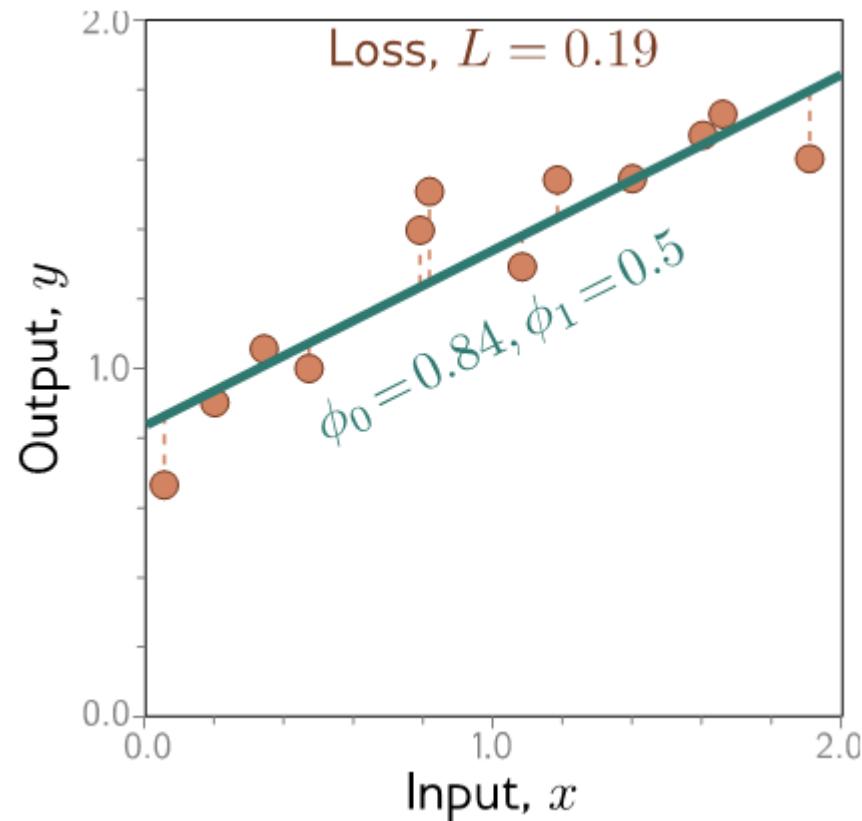
# 2.2 Gradient Descent (GD)

1. Initialize  $\theta$
2. Repeat:
  - Compute gradient  $\nabla_{\theta}\mathcal{L}(\theta)$
  - Update:  $\theta \leftarrow \theta - \eta \nabla_{\theta}\mathcal{L}(\theta)$
3. Until convergence

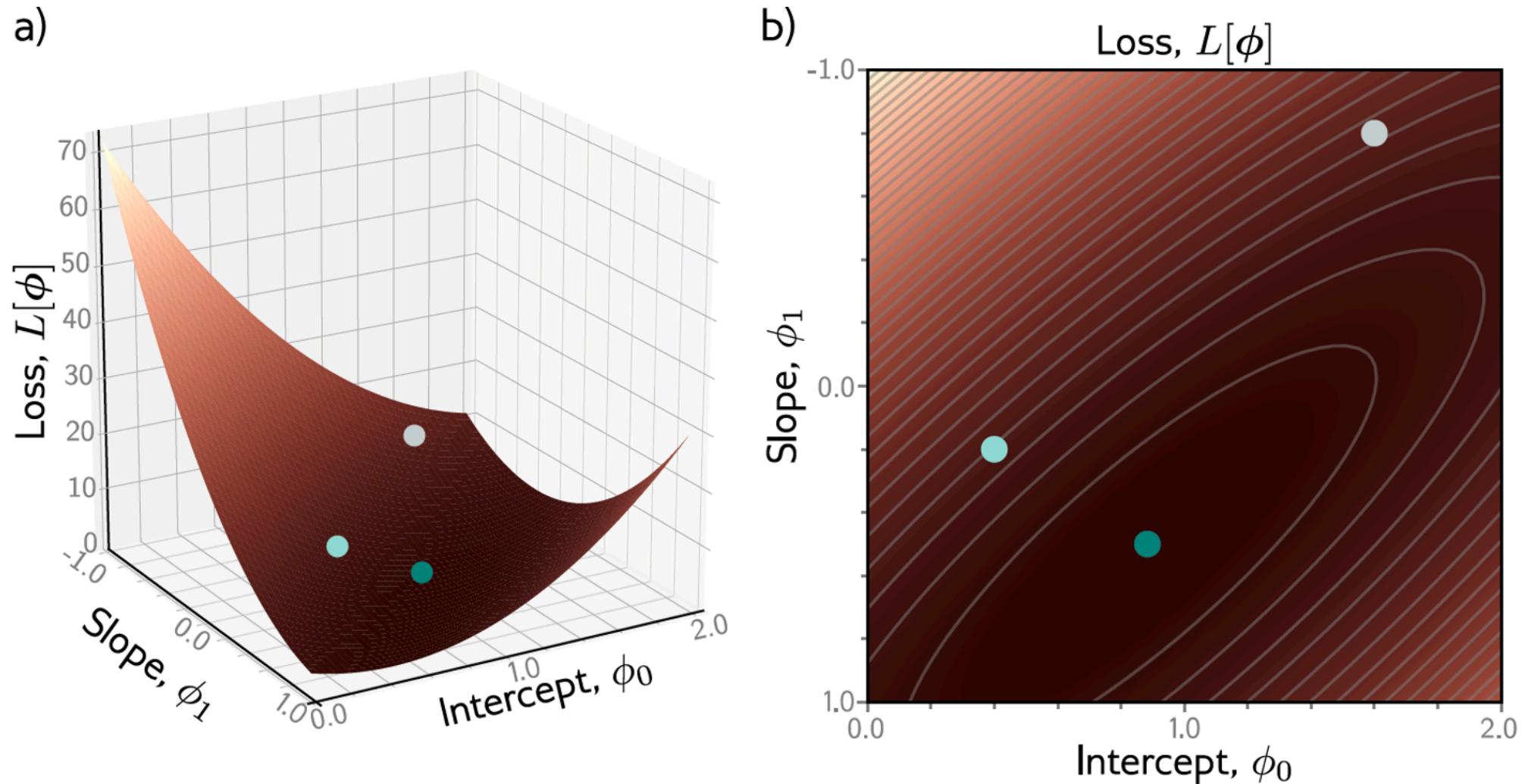
# 2.3 Linear Regression (LR)

1. Model:  $h_\theta(x) = \theta^T x$
2. Loss function:  $\mathcal{L}(h_\theta) = \frac{1}{2} \sum_{i=1}^N (y_i - h_\theta(x_i))^2$
3. Optimization: Find  $\theta^*$  that minimizes  $\mathcal{L}(h_\theta)$ 
  - Gradient descent: Iteratively update  $\theta$  using the gradient of the loss function

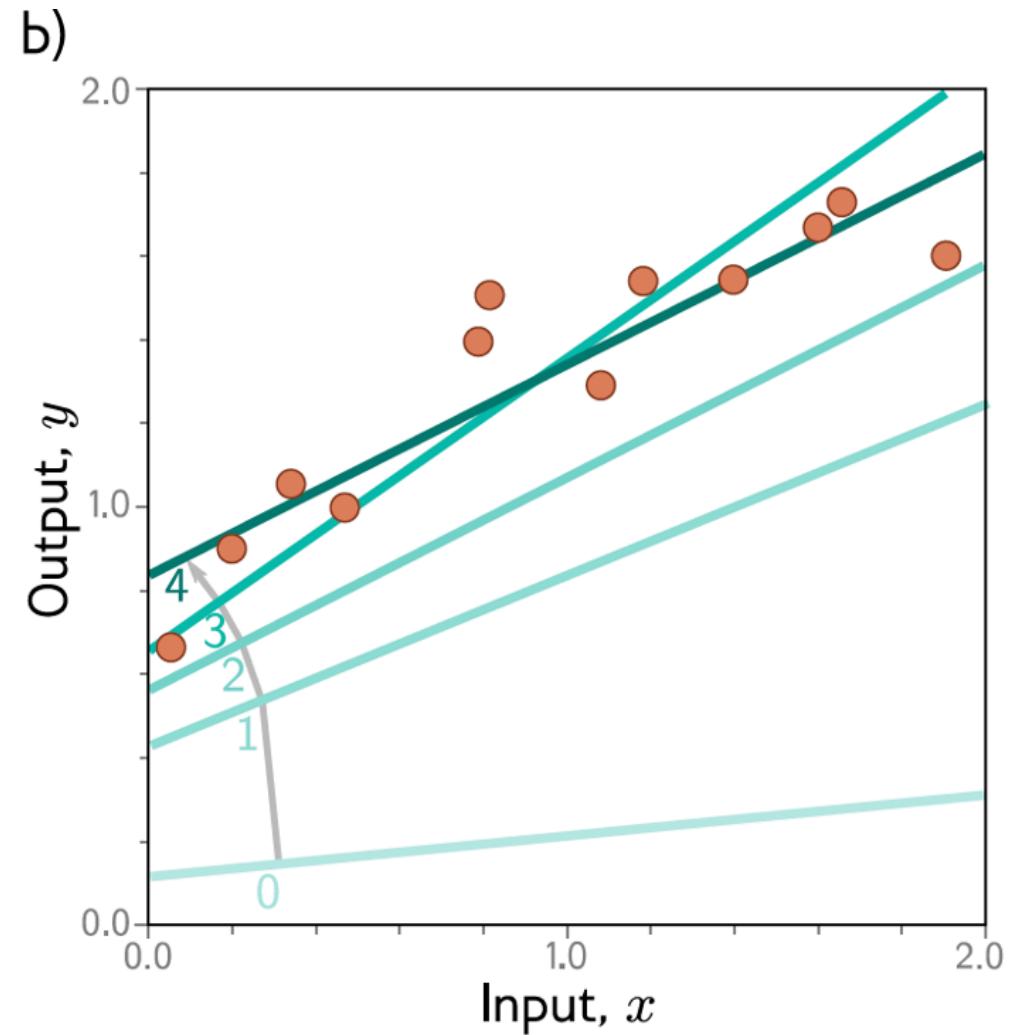
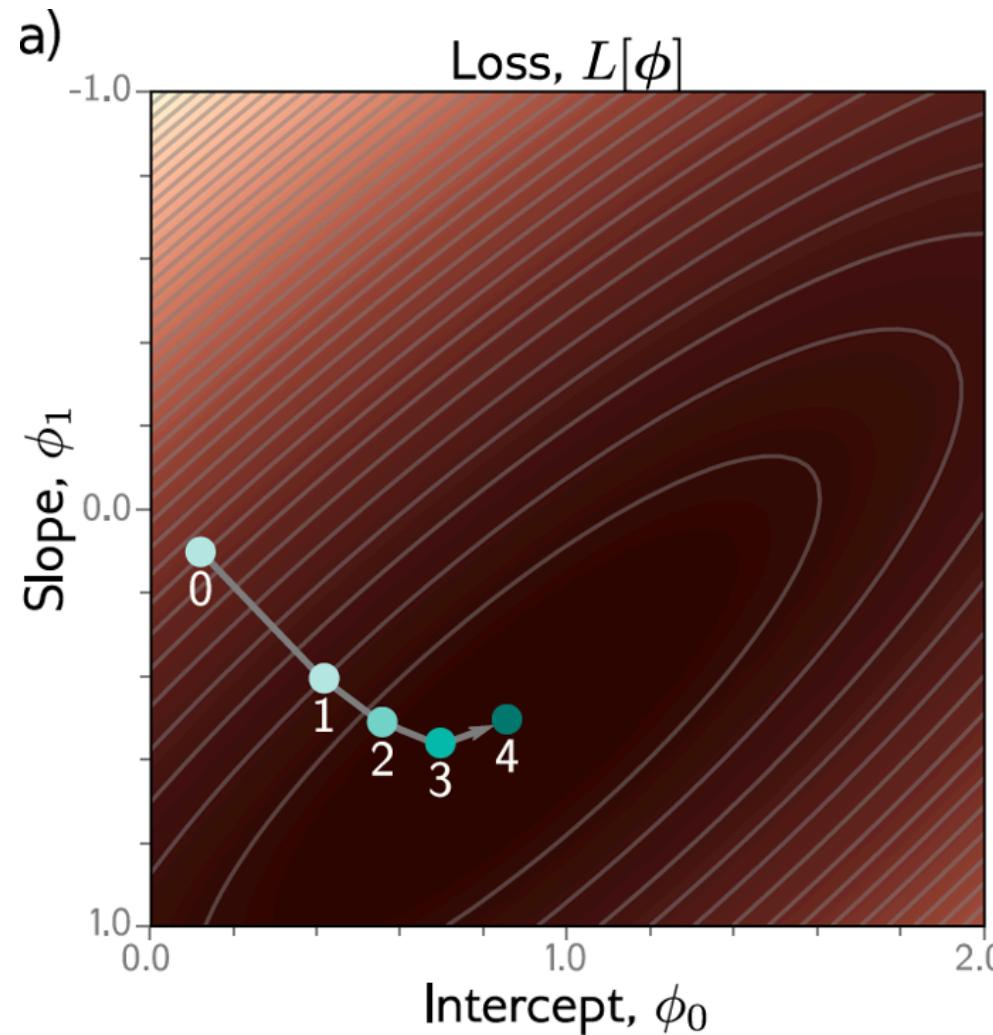
# 2.4 Example: 1D LR loss function



## 2.4 Example: 1D LR loss function

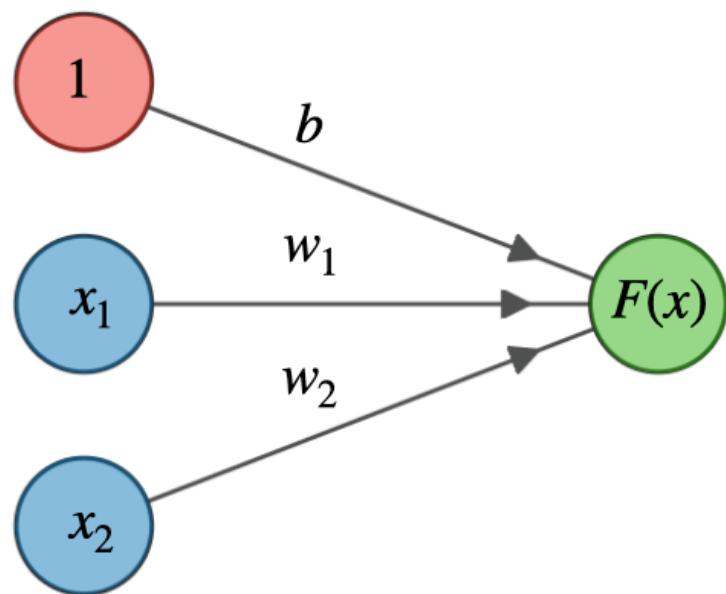


## 2.4 Example: 1D LR loss function



# 2.4 The perceptron: 1D LR as a shallow network

## Linear regression



$$F(x) = w^T x + b$$

Input Layer  $\in \mathbb{R}^3$

Output Layer  $\in \mathbb{R}^1$

$$F(x) = w_1 \cdot x_1 + w_2 \cdot x_2 + 1 \cdot b$$

# 2.5 Possible problems with linear models

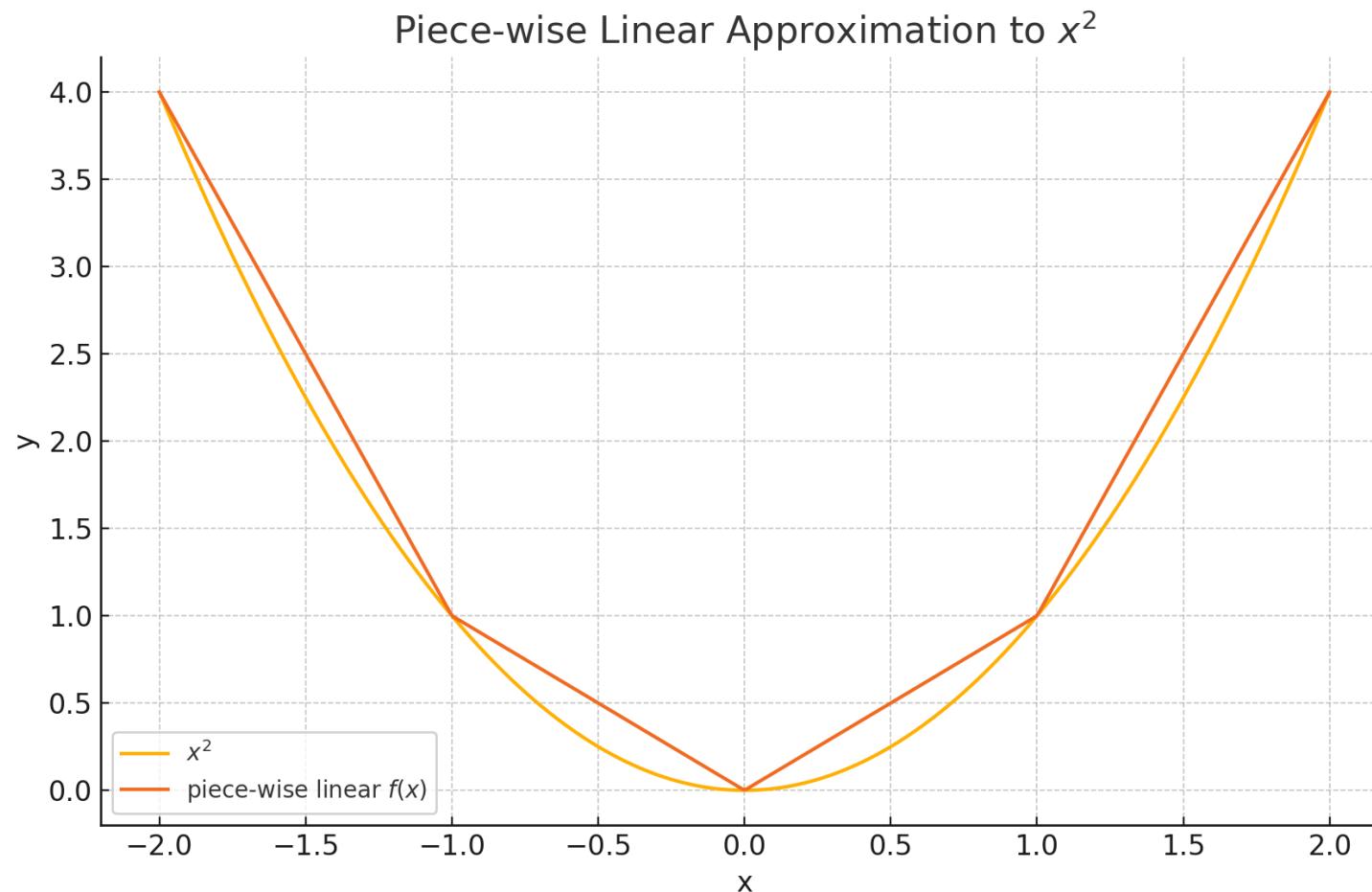
1. Want to be able to describe non-linear relationships
2. Want multiple inputs
3. Want multiple outputs

# 3 Perceptron → Multi-Layered Perceptron (MLP)

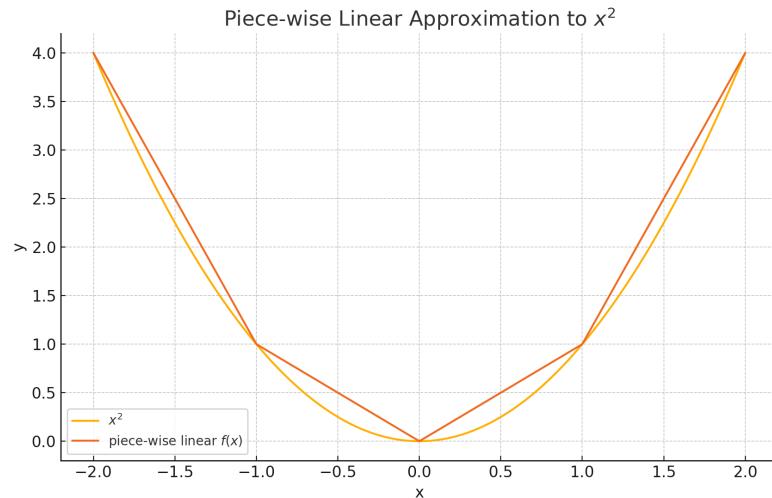
Graphics from 236781 Deep learning course at the Technion  
(Dr. Haim Baskin)

# 3.1 Piecewise linear functions

Consider the following function we would like to learn:



# 3.2 Piecewise linear functions



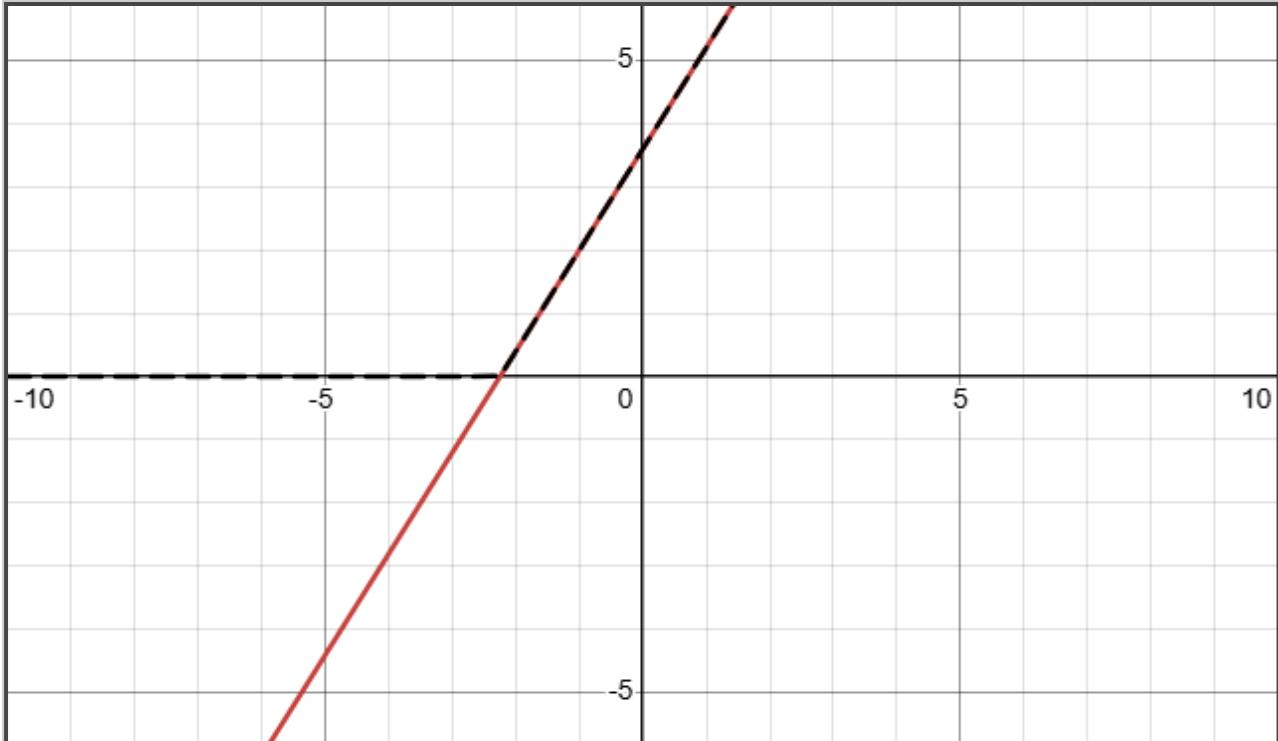
$$f(x) = \begin{cases} -3x - 2, & -2 \leq x \leq -1 \\ -x, & -1 < x \leq 0 \\ x, & 0 < x \leq 1 \\ 3x - 2, & 1 < x \leq 2 \end{cases}$$

- Combinatorial
  - Number of “joints” in the function: 4
  - Location of joints:  $-1, 0, 1$
- Doesn't scale well to higher dimensions
- Differentiability?

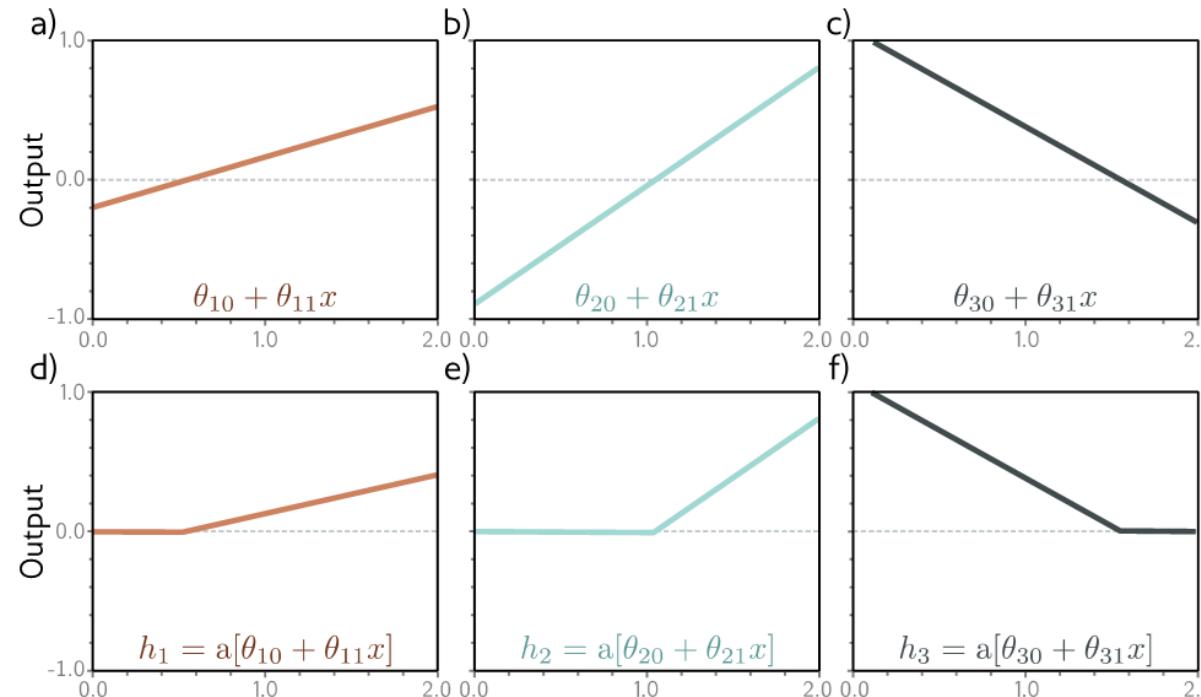
# 3.3 Activation functions: learnable joints!

ReLU =  $\max(0, f(x))$  (Rectified Linear Unit)

$$\text{ReLU}(f(x)) = \begin{cases} 0, & f(x) < 0 \\ f(x), & f(x) \geq 0 \end{cases}$$

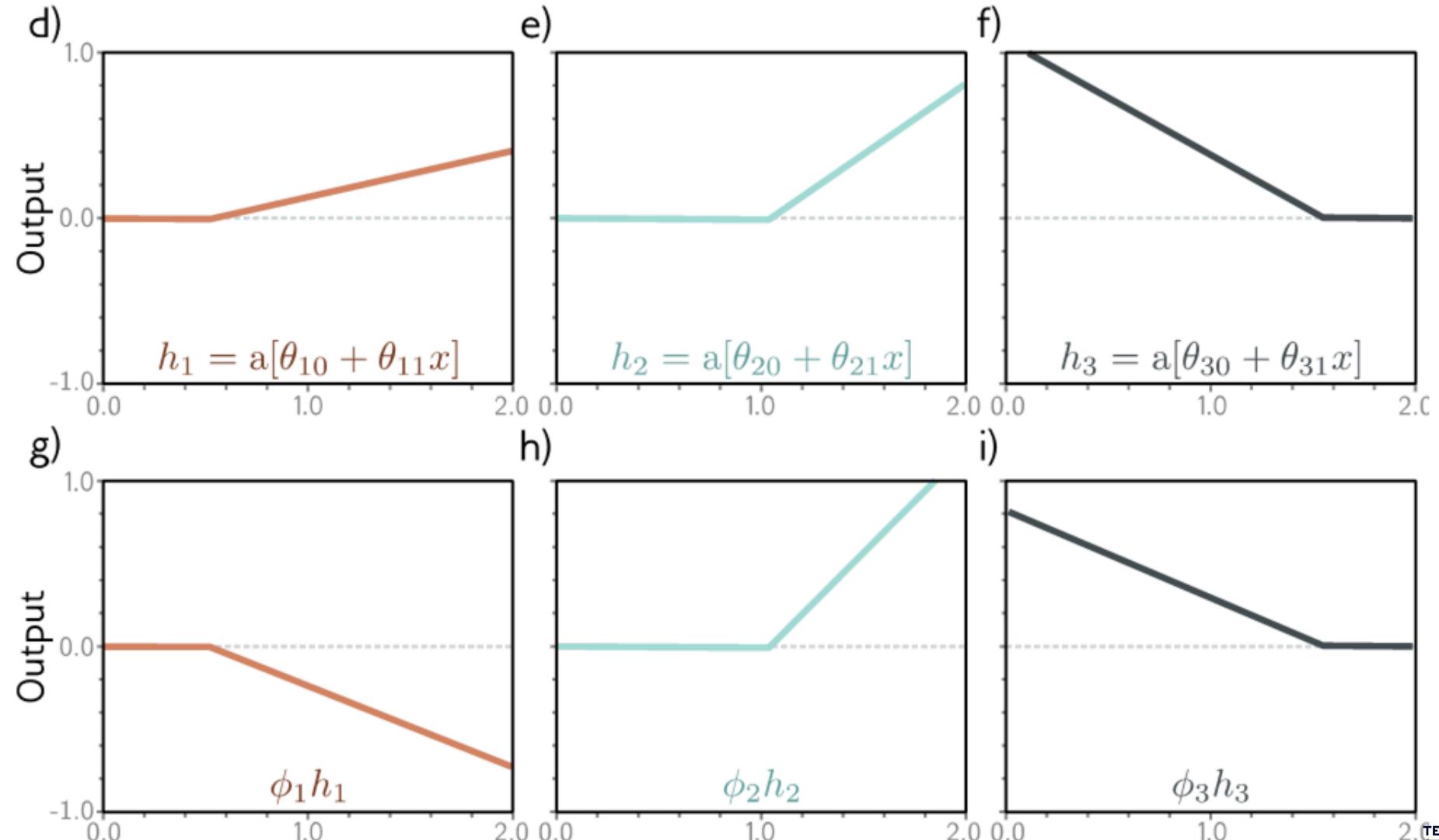


# 3.4 Pass through ReLU functions: hidden units

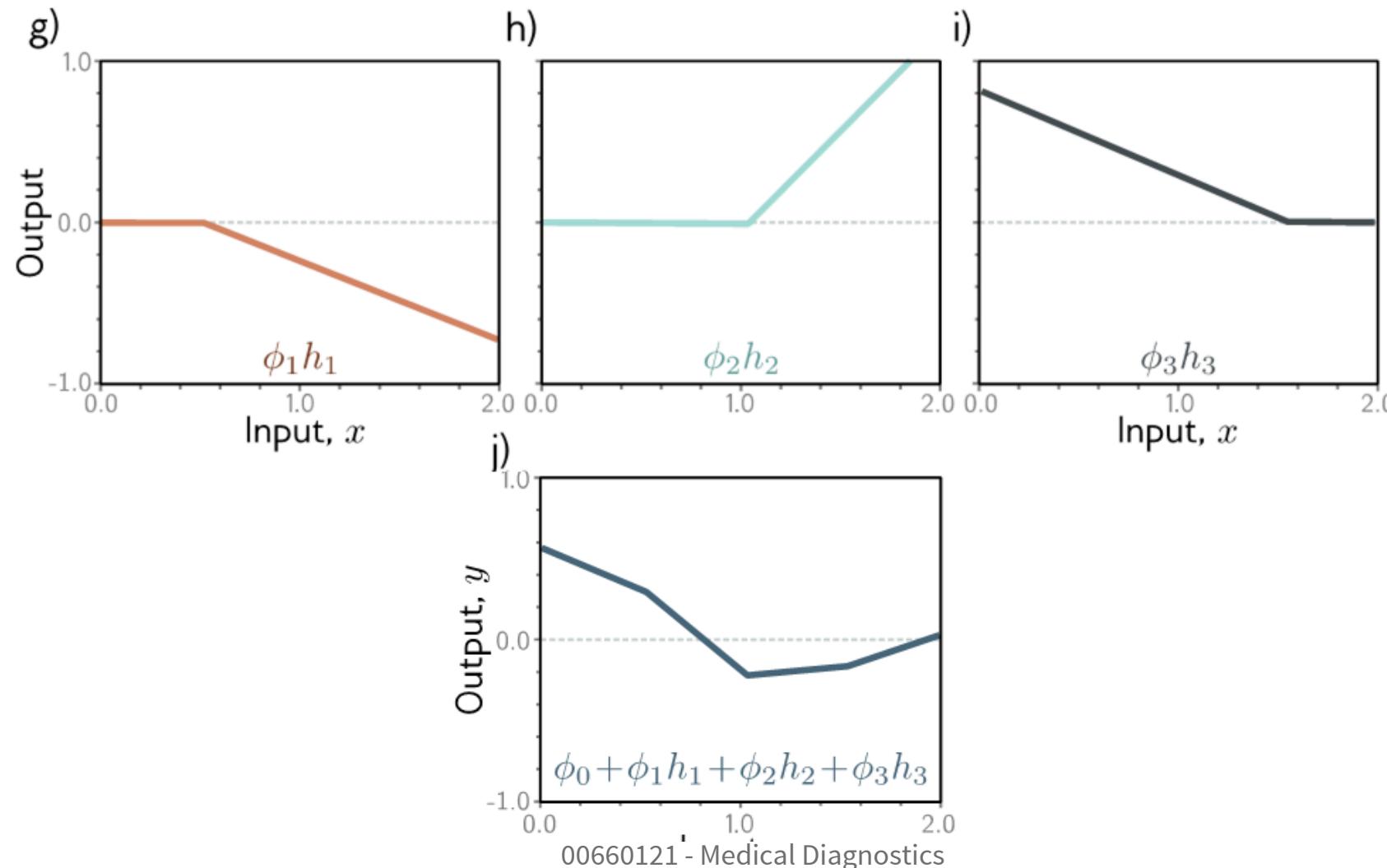


- Other activation functions: sigmoid, tanh, softmax, etc.

# 3.5 Weight the hidden units

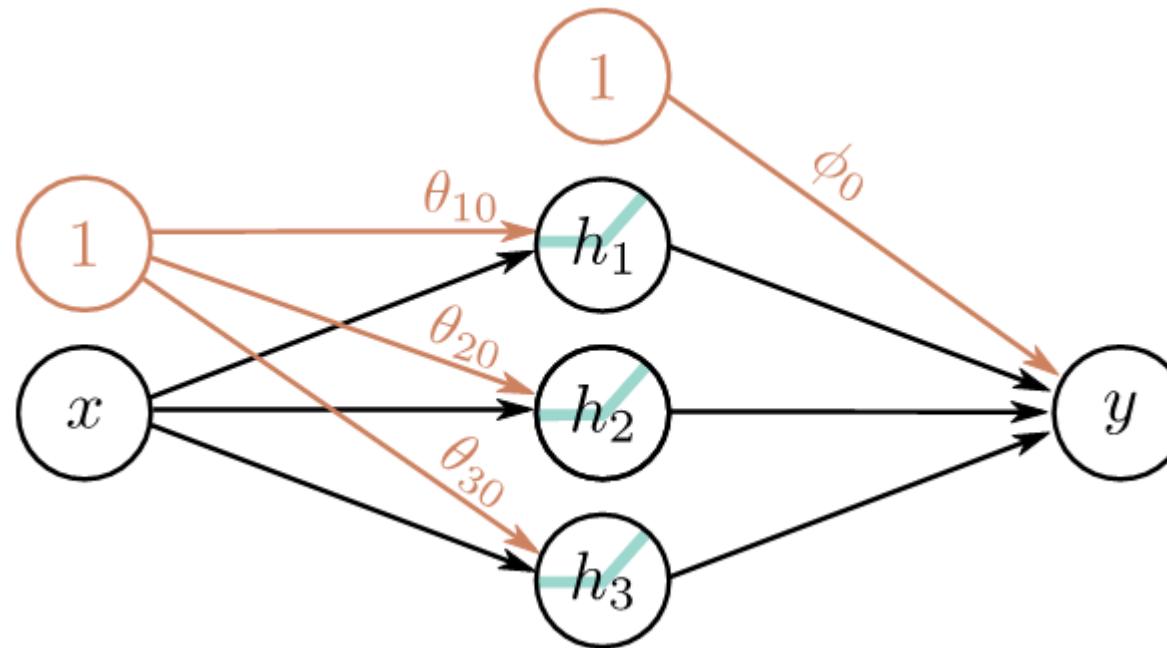


# 3.6 Sum the weighted hidden units to output



$$\begin{aligned}y &= \phi_0 + \phi_1 \operatorname{ReLU}(\theta_{10} + \theta_{11}x) \\&\quad + \phi_2 \operatorname{ReLU}(\theta_{20} + \theta_{21}x) \\&\quad + \phi_3 \operatorname{ReLU}(\theta_{30} + \theta_{31}x)\end{aligned}$$

# 3.7 Lets draw the network again



$$\begin{aligned}
 y = & \phi_0 + \phi_1 \text{ReLU}(\theta_{10} + \theta_{11}x) \\
 & + \phi_2 \text{ReLU}(\theta_{20} + \theta_{21}x) \\
 & + \phi_3 \text{ReLU}(\theta_{30} + \theta_{31}x)
 \end{aligned}$$

00660121 - Medical Diagnostics



# 3.8 Multiple inputs

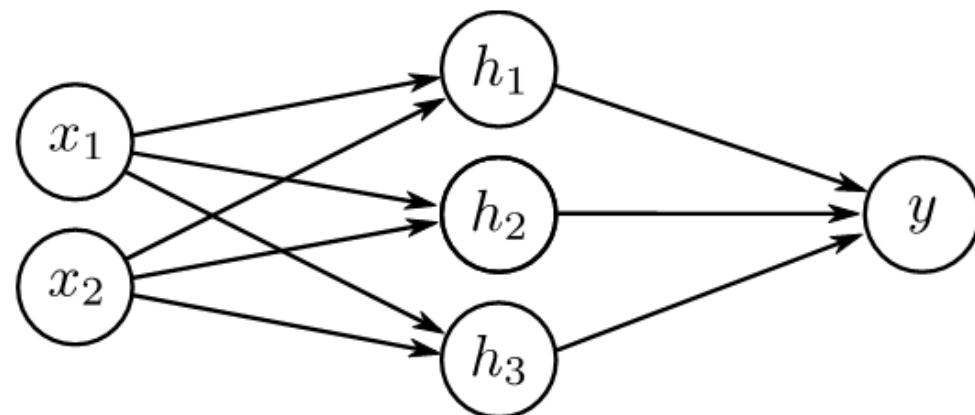
- 2 inputs, 3 hidden units, 1 output

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

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

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

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



# 3.9 Multiple outputs

- 1 input, 4 hidden units, 2 outputs

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

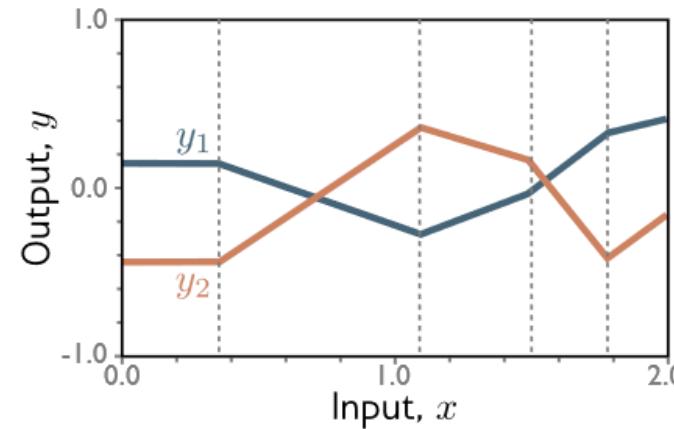
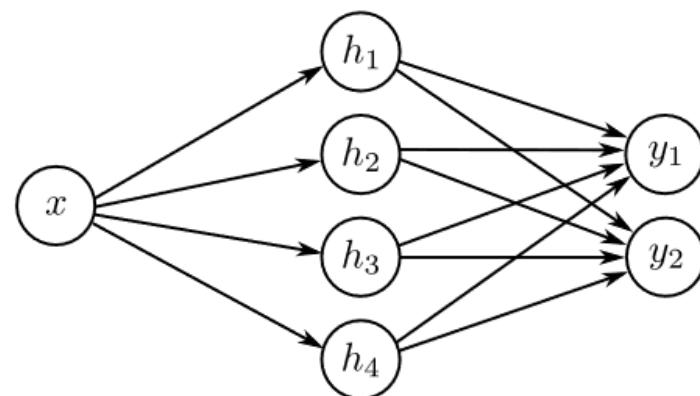
$$h_2 = a[\theta_{20} + \theta_{21}x]$$

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

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

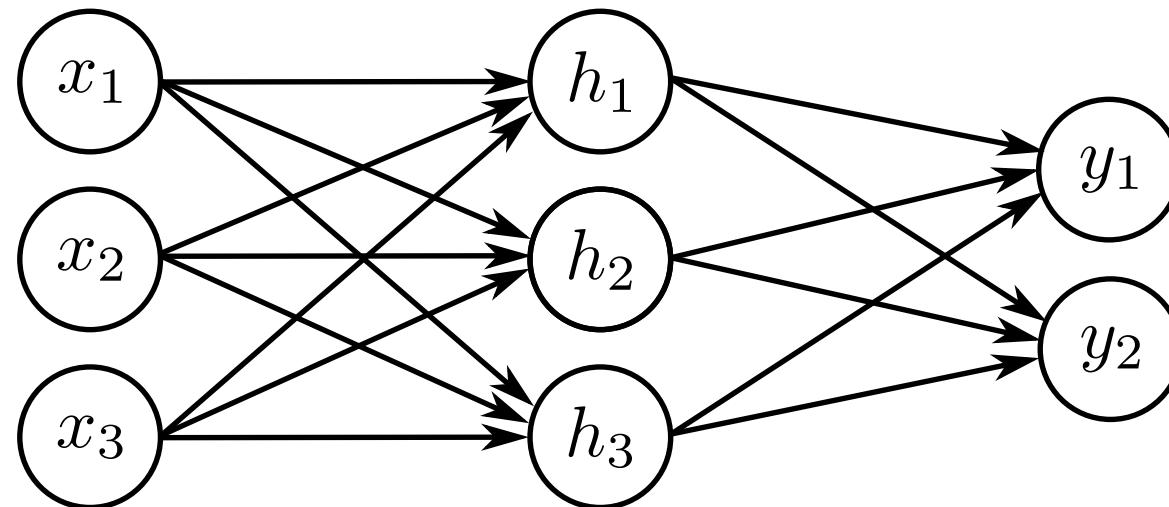
$$y_1 = \phi_{10} + \phi_{11}h_1 + \phi_{12}h_2 + \phi_{13}h_3 + \phi_{14}h_4$$

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

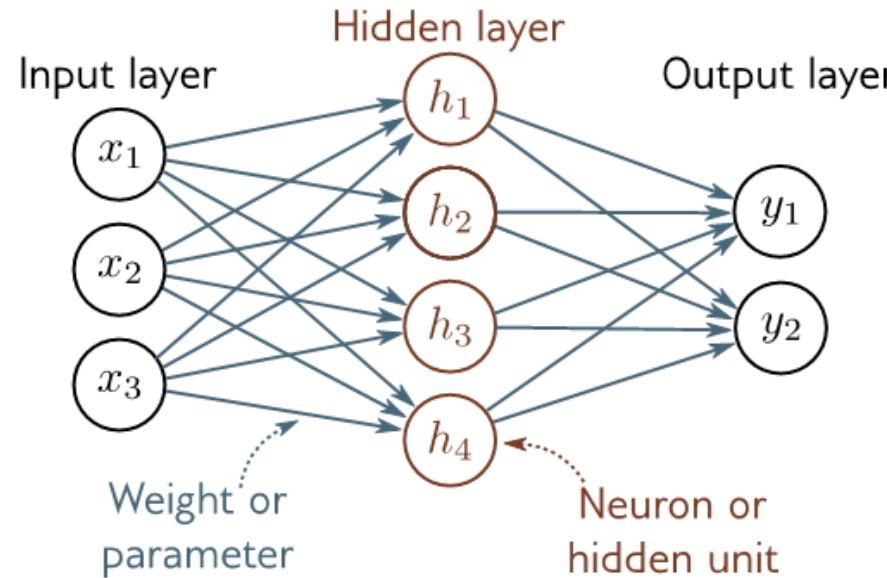


# 3.10 Question

- How many parameters does this model have?



# 3.11 Nomenclature



- Everything in one layer connected to everything in the next = fully connected network
- No loops = feedforward network
- Values after ReLU (activation functions) = activations
- One hidden layer = shallow neural network

# 3.12 MLP: Deep networks

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

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

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

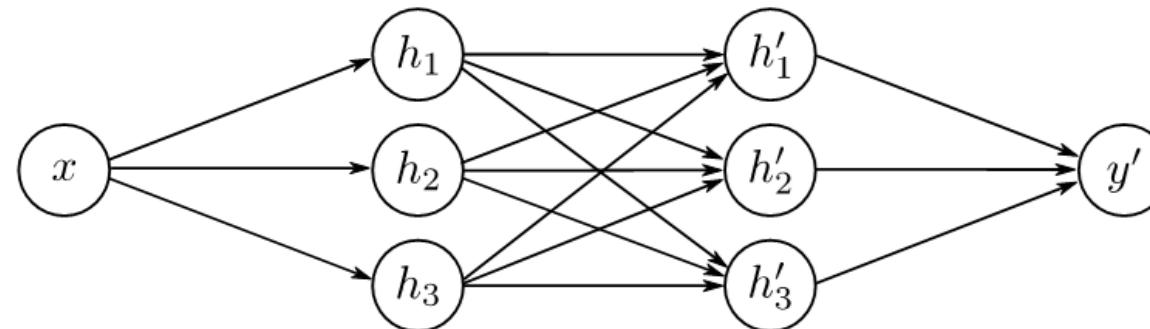
$$h'_1 = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

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

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

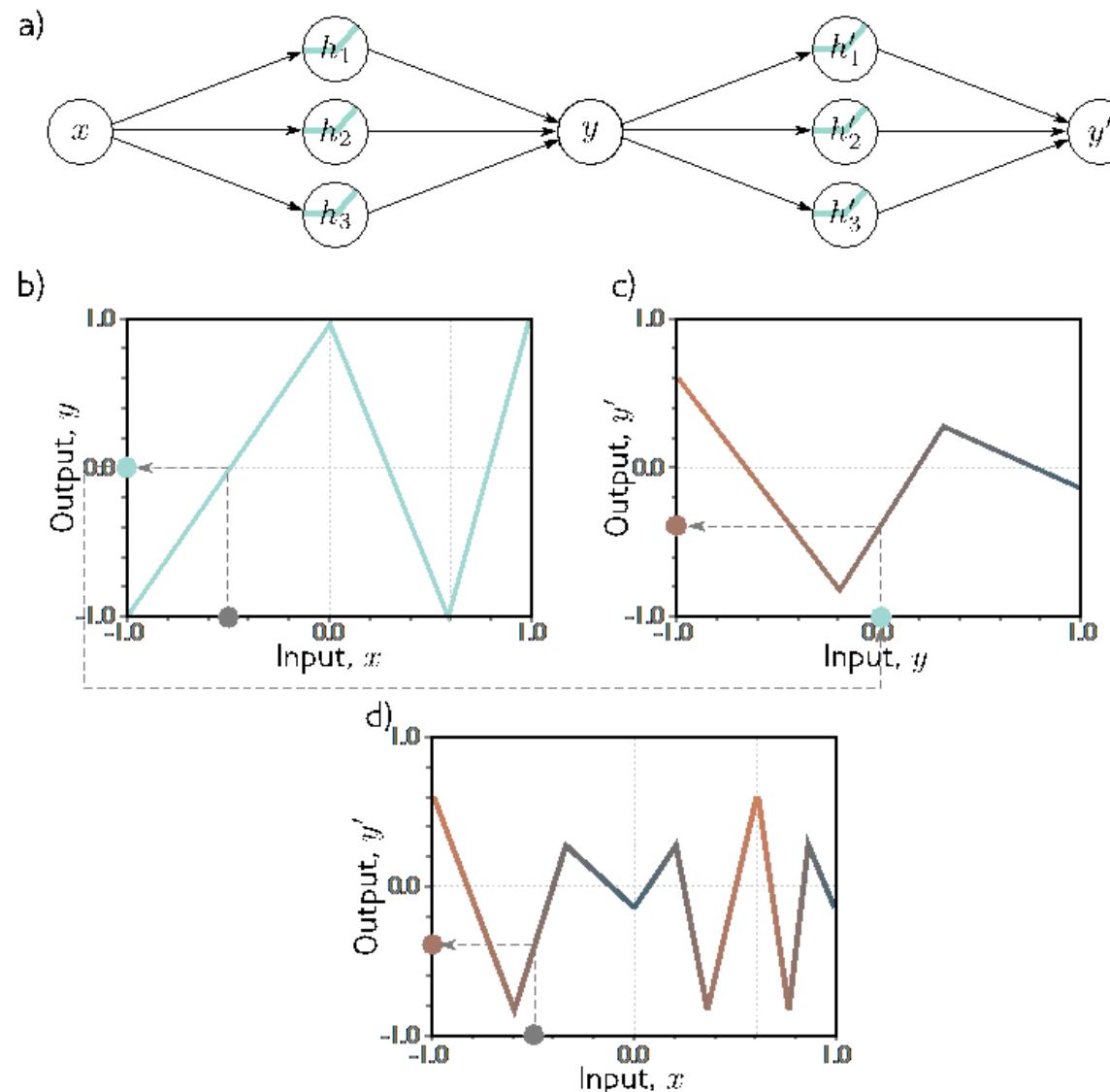
$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$


---

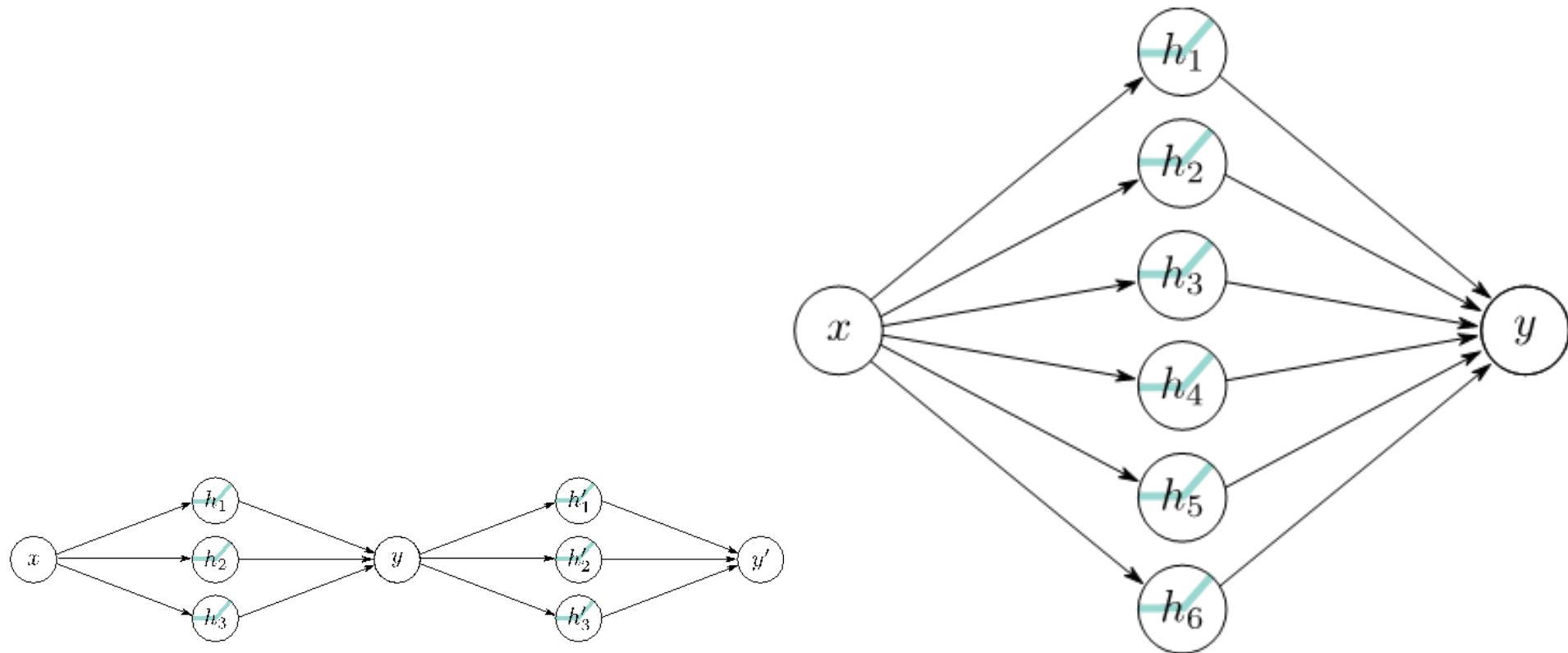


- Hyperparameter: number of hidden layers (Depth)
- Hyperparameter: number of hidden units per layer (Width)

# 3.13 Intuition: Folding

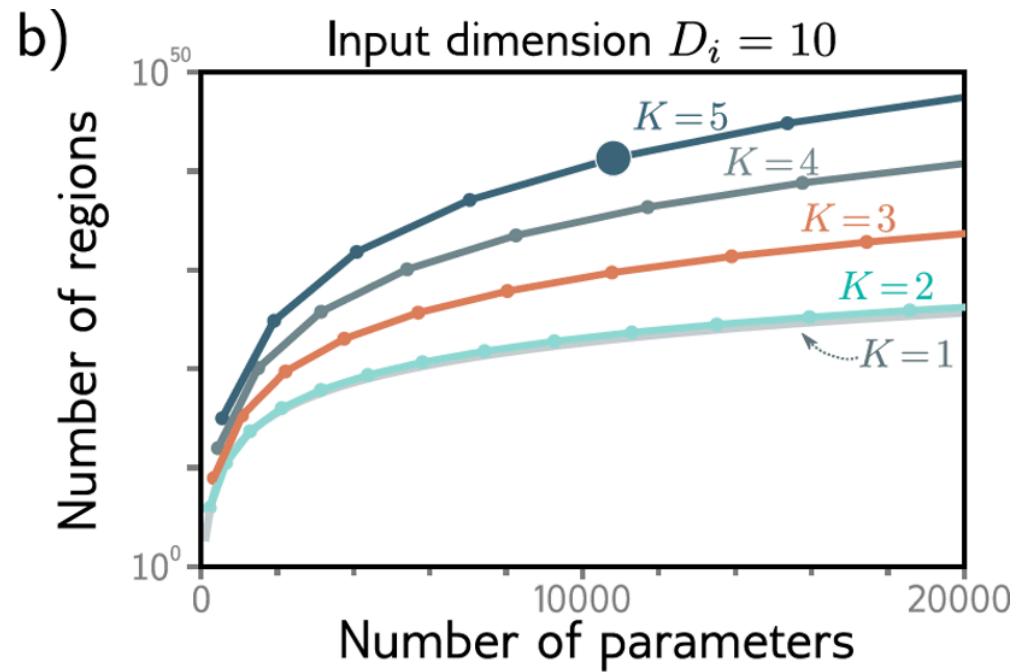
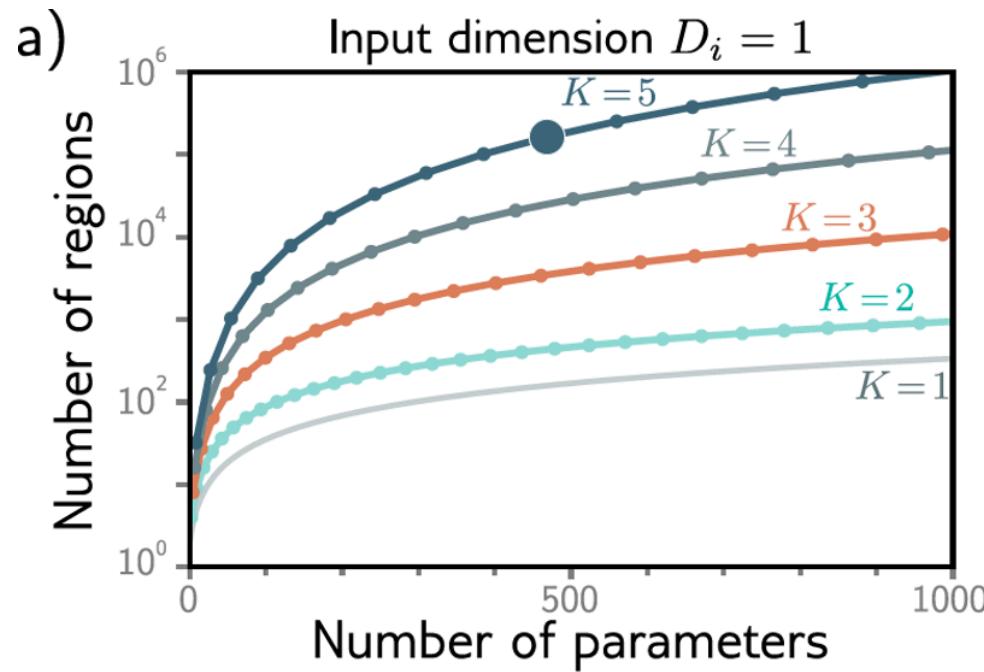


# 3.14 Shallow VS Deep



- Shallow: 19 parameters, max 7 regions
- Deep: 20 parameters, at least 9 regions

# 3.15 Shallow VS Deep visualized



## 3.16 Notations

$$\mathbf{h} = \mathbf{a} \left( \begin{bmatrix} \theta_{10} \\ \theta_{20} \\ \theta_{30} \end{bmatrix} + \begin{bmatrix} \theta_{11} \\ \theta_{21} \\ \theta_{31} \end{bmatrix} x \right) = \mathbf{a}(\boldsymbol{\theta}_0 + \boldsymbol{\theta}x)$$

$$\mathbf{h}' = \mathbf{a} \left( \begin{bmatrix} \psi_{10} \\ \psi_{20} \\ \psi_{30} \end{bmatrix} + \begin{bmatrix} \psi_{11} & \psi_{12} & \psi_{13} \\ \psi_{21} & \psi_{22} & \psi_{23} \\ \psi_{31} & \psi_{32} & \psi_{33} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} \right) = \mathbf{a}(\boldsymbol{\psi}_0 + \boldsymbol{\psi} \mathbf{h})$$

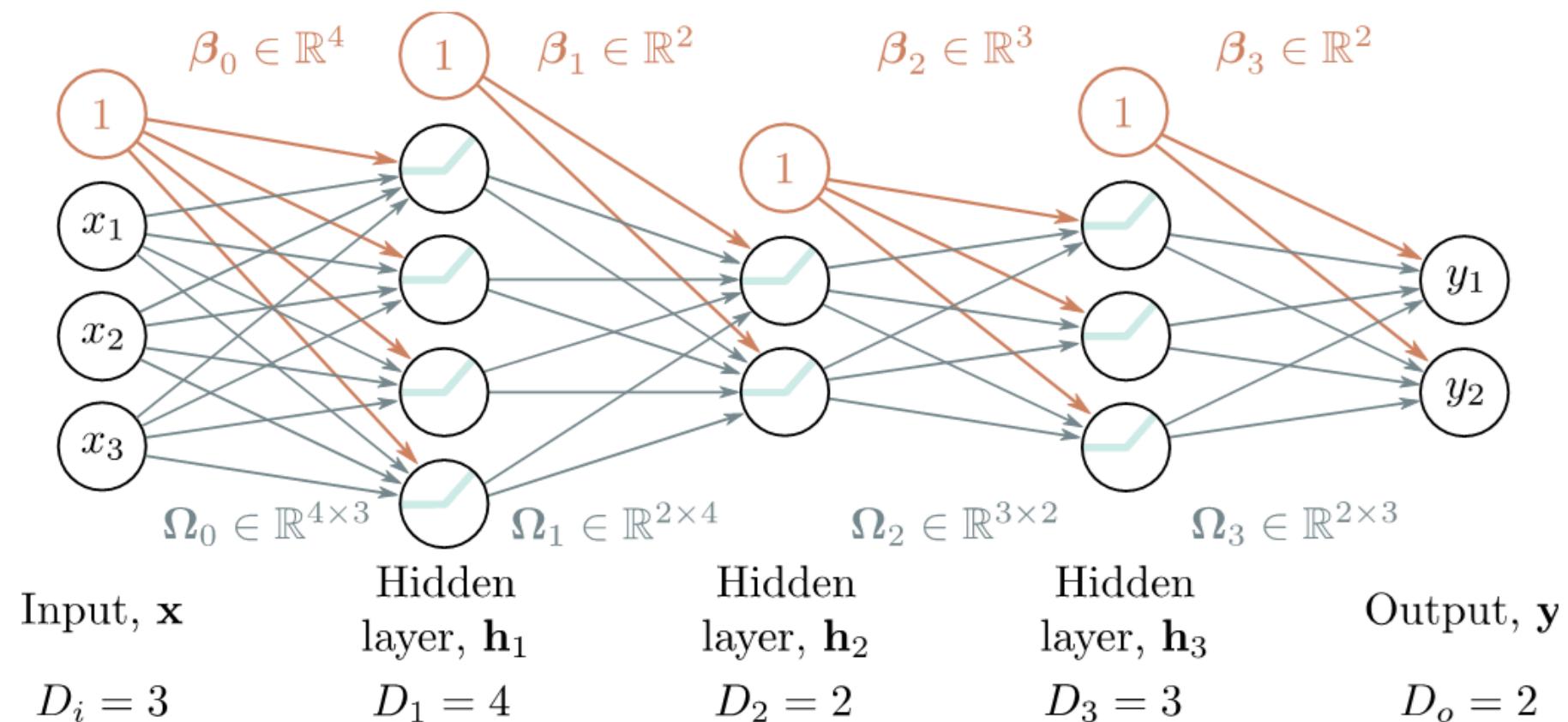
$$y = [\phi_1 \quad \phi_2 \quad \phi_3] \begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix} + \phi_0 = \boldsymbol{\phi}^\top \mathbf{h}' + \phi_0$$

00660121 - Medical Diagnostics

00660121 - Medical Diagnostics

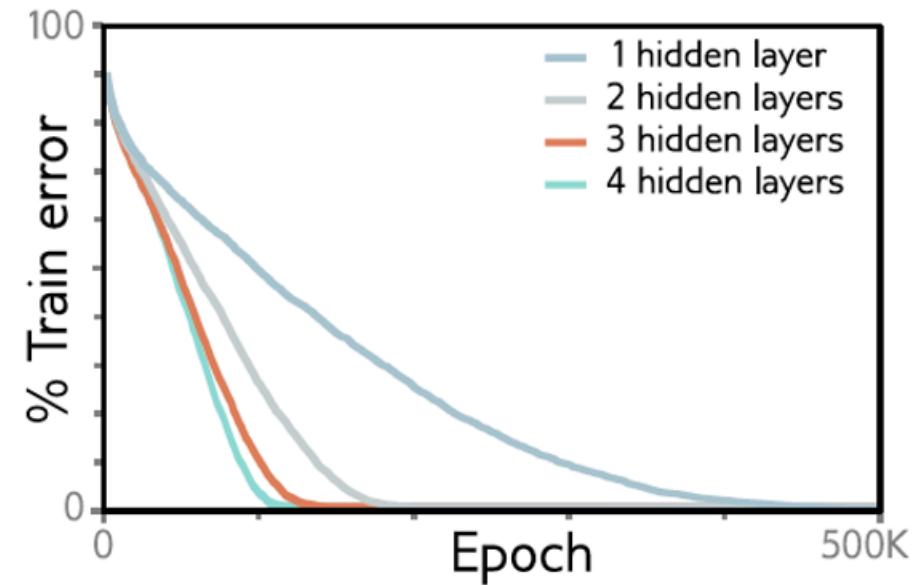


# 3.17 Deep networks



# 3.18 Deep networks

**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.



# 3.19 Deep Neural Network playground

00660121 - Medical Diagnostics



# 4 Backpropagation

How do we find  $\nabla \mathcal{L}(h_\theta)$  for deep networks?

$$L'(\theta, \psi, \phi) = \dots ?$$

# 4.1 Chain rule (1-D recap)

If  $y = f(g(x))$   
and  $u = g(x)$  then

$$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx}.$$

*One derivative “chain” per link.*

# 4.2 Chain rule example

Let  $u = g(x) = x^2$ ,  $y = f(u) = \sin u$ ,  $x = 0.3$

- In Calculus class:
  - $y = f(g(x)) = \sin(x^2)$
  - $f'(x) = \cos(x^2) \cdot 2x = \cos(u) \cdot u'$
- Forward
  - $u = g(0.3) = 0.09$ ,  $y = f(u) = \sin 0.09 \approx 0.0899$
- Gradients
  - $\frac{dy}{du} = \cos u = 0.996$  ,  $\frac{du}{dx} = 2x = 0.6$
  - $\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx} = 0.996 \times 0.6 \approx 0.598$

# 4.3 Motivation for backpropagation

Take the toy model:

$$f(x, y, z) = (x + y) z = qz, \quad x, y \xrightarrow{+} q \xrightarrow{\times z} f.$$

To get every gradient we build three separate chains, what is the problem?

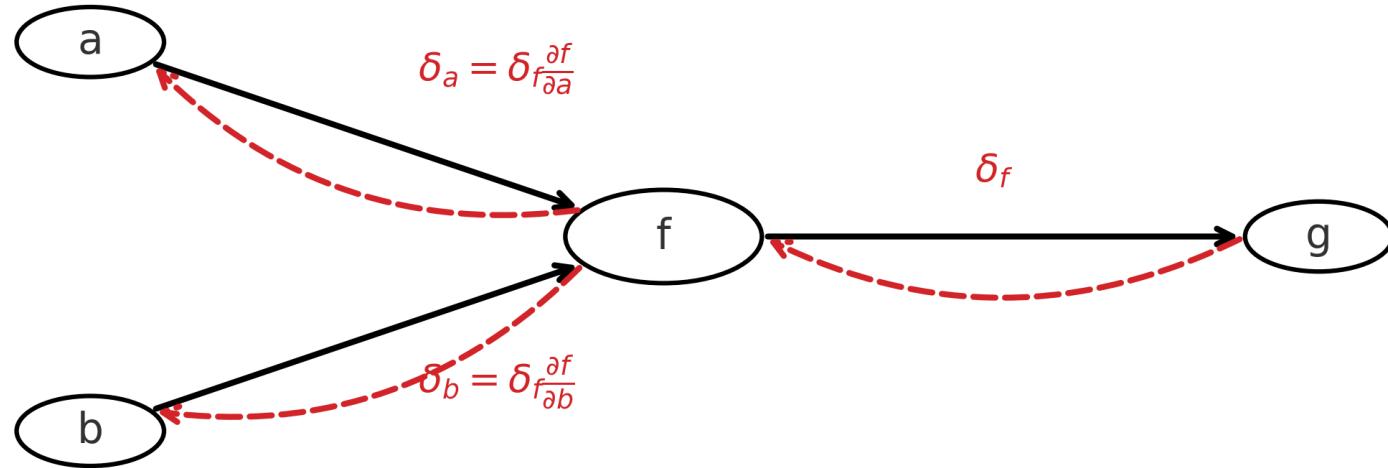
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x} = z \cdot 1 = z$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y} = z \cdot 1 = z$$

$$\frac{\partial f}{\partial z} = \frac{\partial f}{\partial z} = q$$

# 4.3 Backpropagation

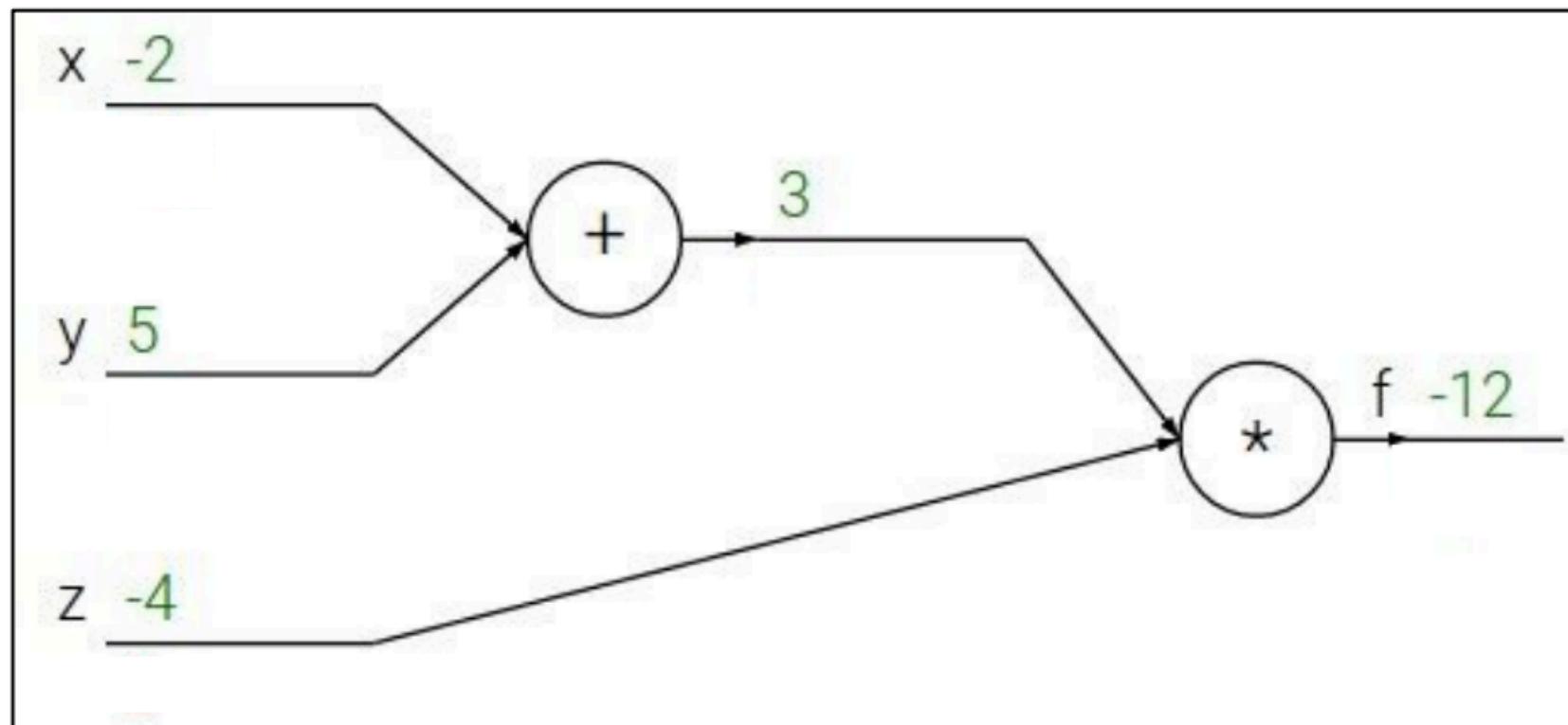
- We need to compute the gradient of the loss function with respect to each parameter in the network.  $\frac{\partial L}{\partial \theta}$ ,  $\frac{\partial L}{\partial \psi}$ ,  $\frac{\partial L}{\partial \phi}$
- Use the chain rule to compute gradients efficiently by reusing intermediate results from the forward pass.



## 4.4 Backpropagation example

$$f(x, y, z) = (x + y)z = qz, \quad x, y \xrightarrow{+} q \xrightarrow{\times z} f.$$

- Let  $x = -2, y = 5, z = -4$ . Need



# 4.5 Backprop algorithm (simplified)

1. Build the computational graph
2. Forward pass: compute all intermediate values
3. Backward pass: compute gradients (once!) using the chain rule
4. Update parameters using gradients

# 5 Implementing and training models

# 5.1 Defining the model

```
1 W1 = HeNormal([h1_dim, in_dim])      # weight matrix
2 b1 = zeros([h1_dim, 1])
3
4 W2 = HeNormal([h2_dim, h1_dim])
5 b2 = zeros([h2_dim, 1])
6
7 W3 = HeNormal([out_dim, h2_dim])
8 b3 = zeros([out_dim, 1])
9
10 forward(X):
11     z1 = W1 @ X + b1
12     a1 = ReLU(z1)
13
14     z2 = W2 @ a1 + b2
15     a2 = ReLU(z2)
16
17     z3 = W3 @ a2 + b3
18     v̂ = z3
```

# 5.2 The deep learning training loop

```
1 import loss_function, SGD, backprop, forward
2
3 optimizer = SGD(learning_rate=0.01)
4
5 for epoch in range(num_epochs):
6     optimizer.zero_grad()
7
8     # Forward pass
9     Y_hat = forward(X)
10    loss = loss_function(Y_hat, Y)
11
12    # Backward pass
13    gradients = backprop(loss, w1, b1, w2, b2, w3, b3)
14
15    # Update parameters
16    optimizer.step(gradients)
```