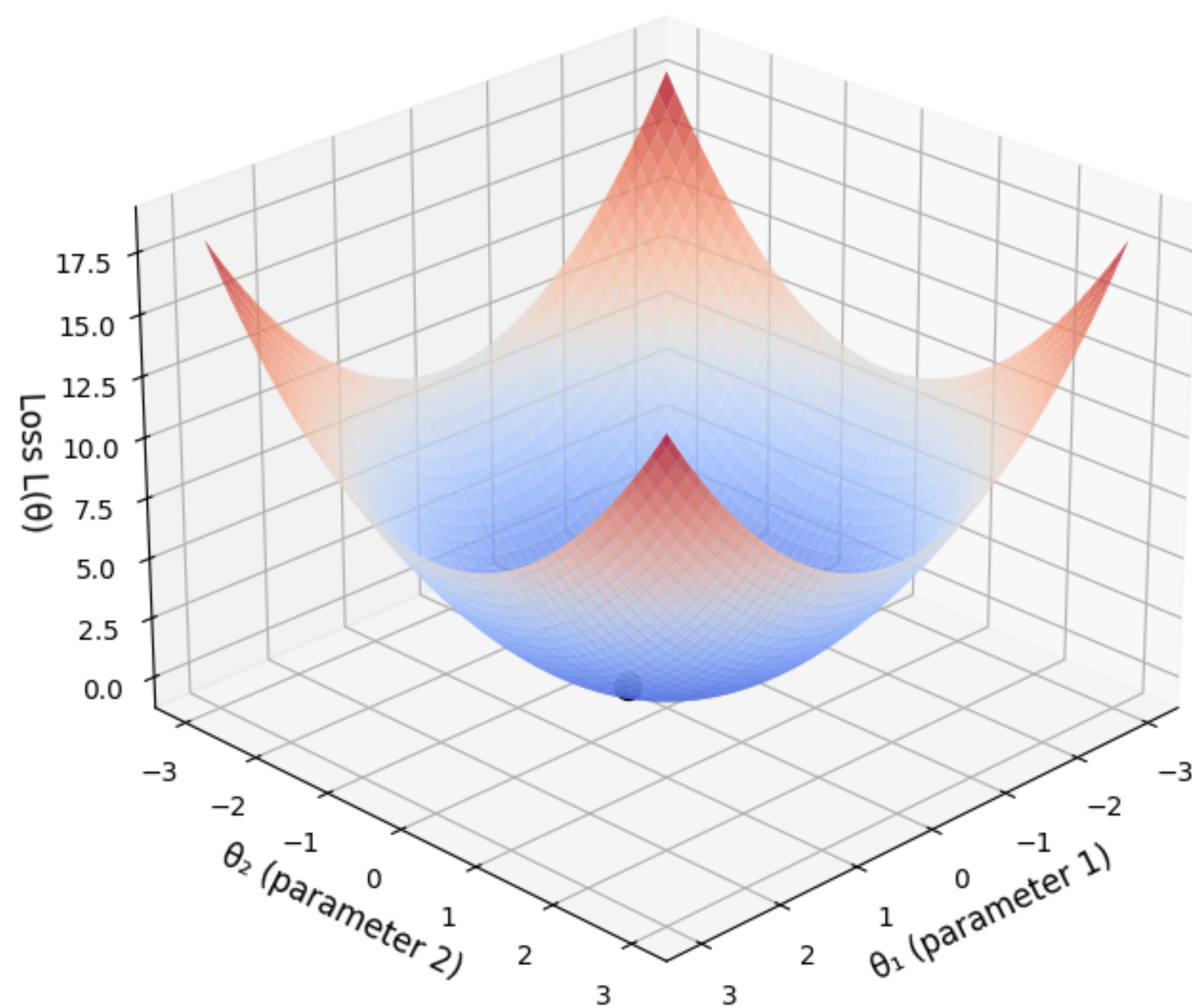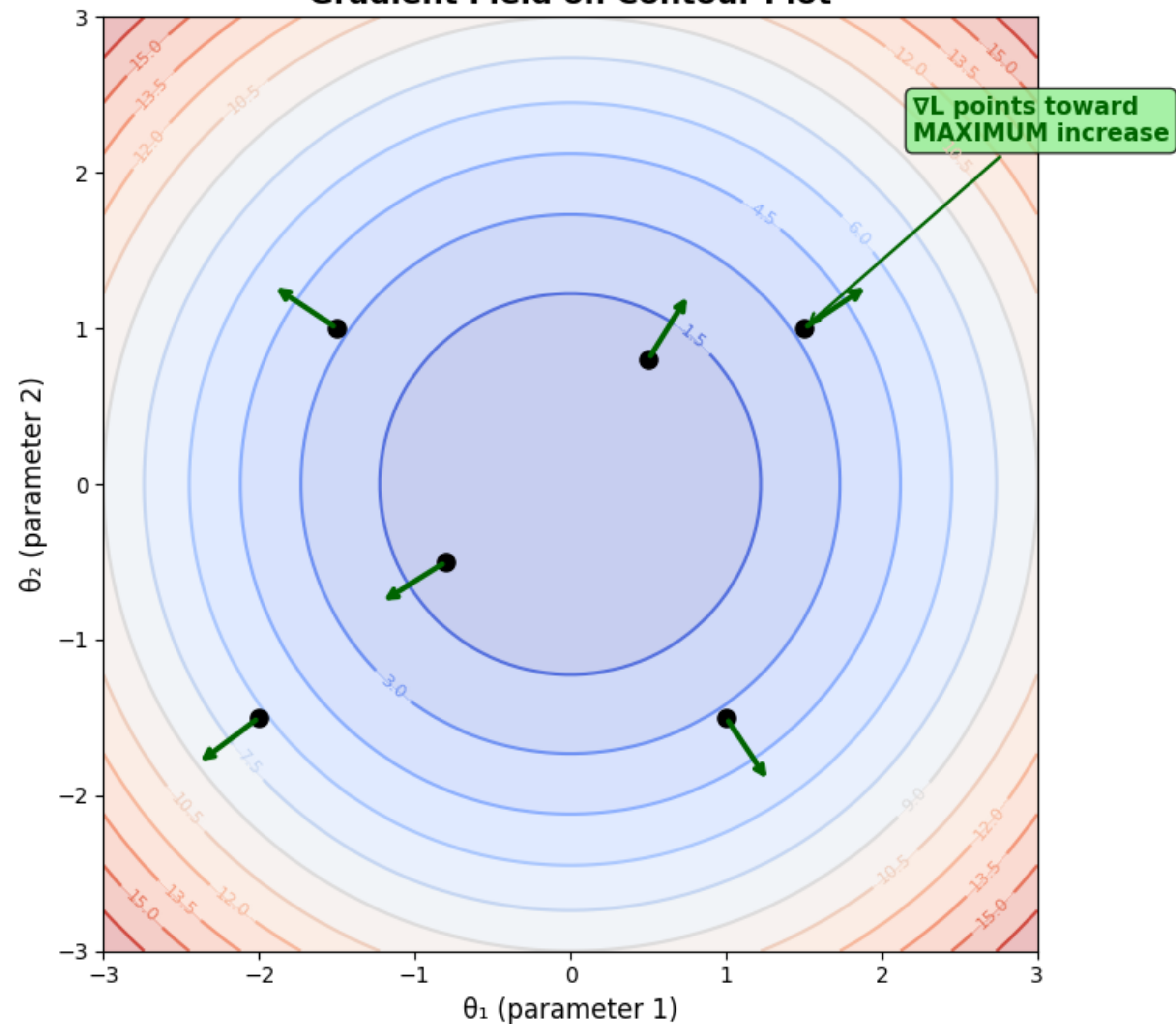# Deep Learning (1470)

**Randall Balestriero**

**Class 4**

# Rewind!

# Gradient: Intuition



3D Loss Surface

Gradient Field on Contour Plot

∇L points toward MAXIMUM increase

# Gradient: Intuition



## Linear Classifier Decision Boundary

Iteration 0

- Class 0
- Class 1

## Training Loss Evolution

Loss: 0.3059

# Gradient: Intuition



**Linear Classifier Decision Boundary**

**Training Loss Evolution**

Loss: 0.3059

# MNIST

The most famous dataset in Deep Learning

**M**odified **N**ational **I**nstitute of **S**tandards and **T**echnology database
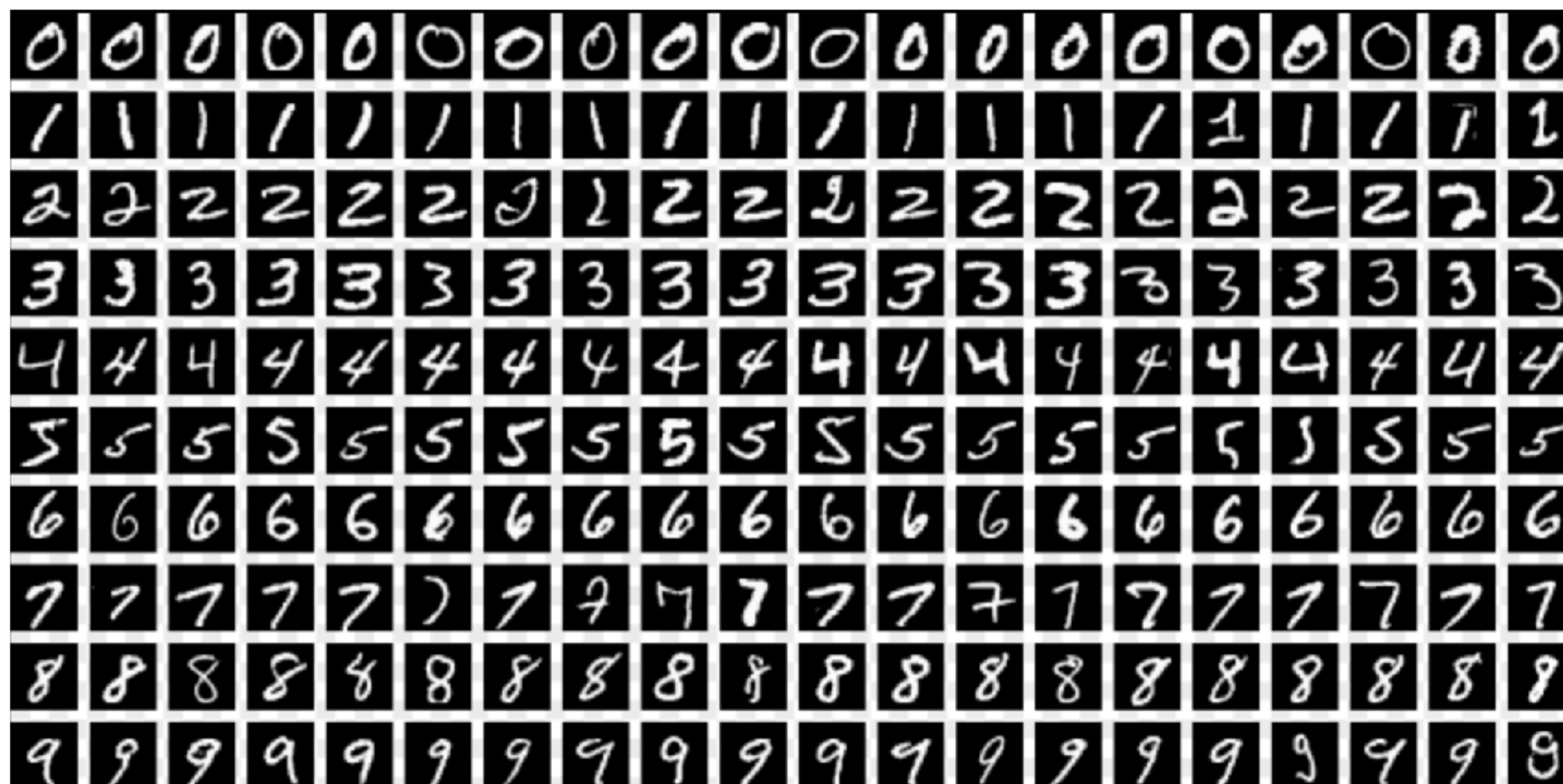


Image courtesy of Wikipedia

# MNIST

# MNIST

- What is D (dimension) of the samples?

# MNIST

- What is D (dimension) of the samples?

- How many classes do we have?

# MNIST

- What is D (dimension) of the samples?

- How many classes do we have?

- Do you think the classes are linearly separable?

# MNIST

- What is D (dimension) of the samples?

- How many classes do we have?

- Do you think the classes are linearly separable?

A linear model reaches 94% accuracy!

# Your First Deep Network!

# Your First Deep Network!

- Why do we want to use a Deep Network?

# Your First Deep Network!

- Why do we want to use a Deep Network?

- Could you think of "features" to extract by-hand that would improve linear model accuracy?

# Your First Deep Network!

- Why do we want to use a Deep Network?

- Could you think of "features" to extract by-hand that would improve linear model accuracy?

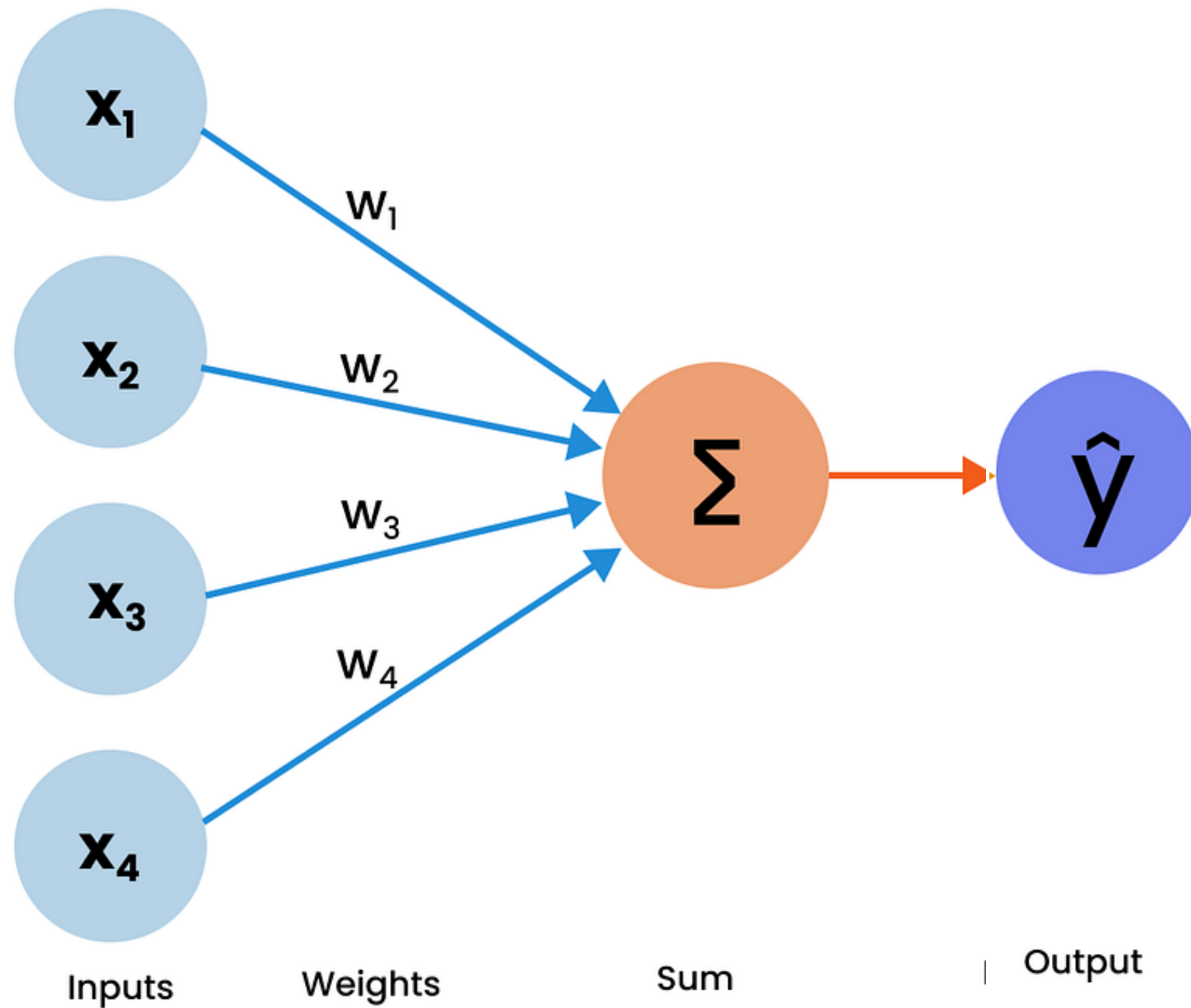- Which accuracy do you think we can reach?

# Your First Deep Network!

- Why do we want to use a Deep Network?

- Could you think of "features" to extract by-hand that would improve linear model accuracy?
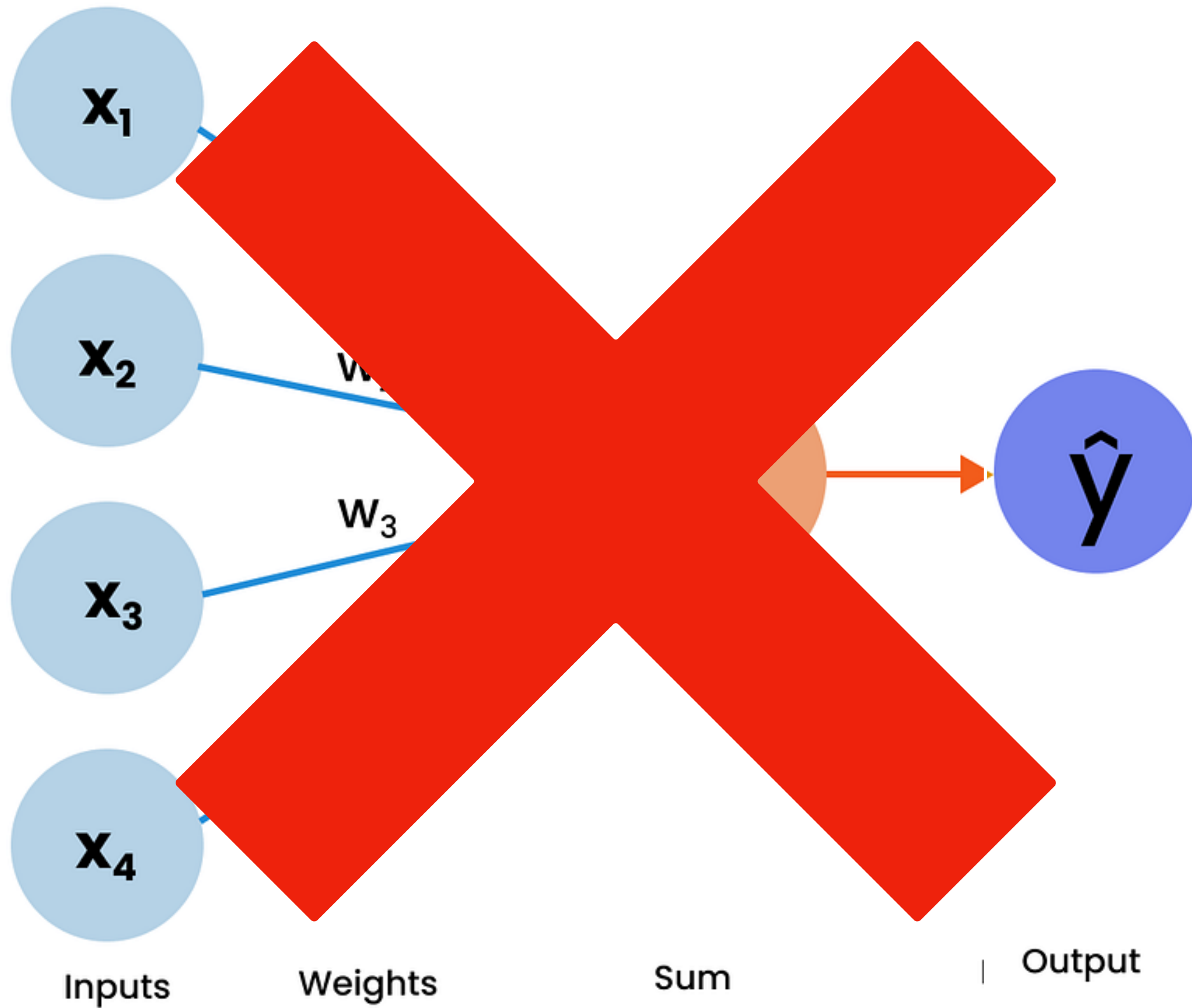
- Which accuracy do you think we can reach?

A Deep Network reaches 99.5% accuracy!

# Your First Deep Network!

# Your First Deep Network!



Inputs      Weights      Sum      Output

# Your First Deep Network!

# Your First Deep Network!



## 3-Layer MLP with ReLU Activations

$$ReLU(z) = max(0, z)$$

| Input Layer (764 units) | Hidden Layer 1 (ReLU) | Hidden Layer 2 (ReLU) | Output Layer (10 units) |

# Your First Deep Network: Why?



3-Layer MLP with ReLU Activations

ReLU(z) = max(0, z)

$W^{(1)}$    $W^{(2)}$    $W^{(3)}$

Input Layer
(764 units)

Hidden Layer 1
(ReLU)

Hidden Layer 2
(ReLU)

Output Layer
(10 units)

# Your First Deep Network: Why?



### 3-Layer MLP with ReLU Activations

$ReLU(z) = max(0, z)$

Input Layer
(764 units)

Hidden Layer 1
(ReLU)

Hidden Layer 2
(ReLU)

Output Layer
(10 units)

- Why 2 hidden layers?
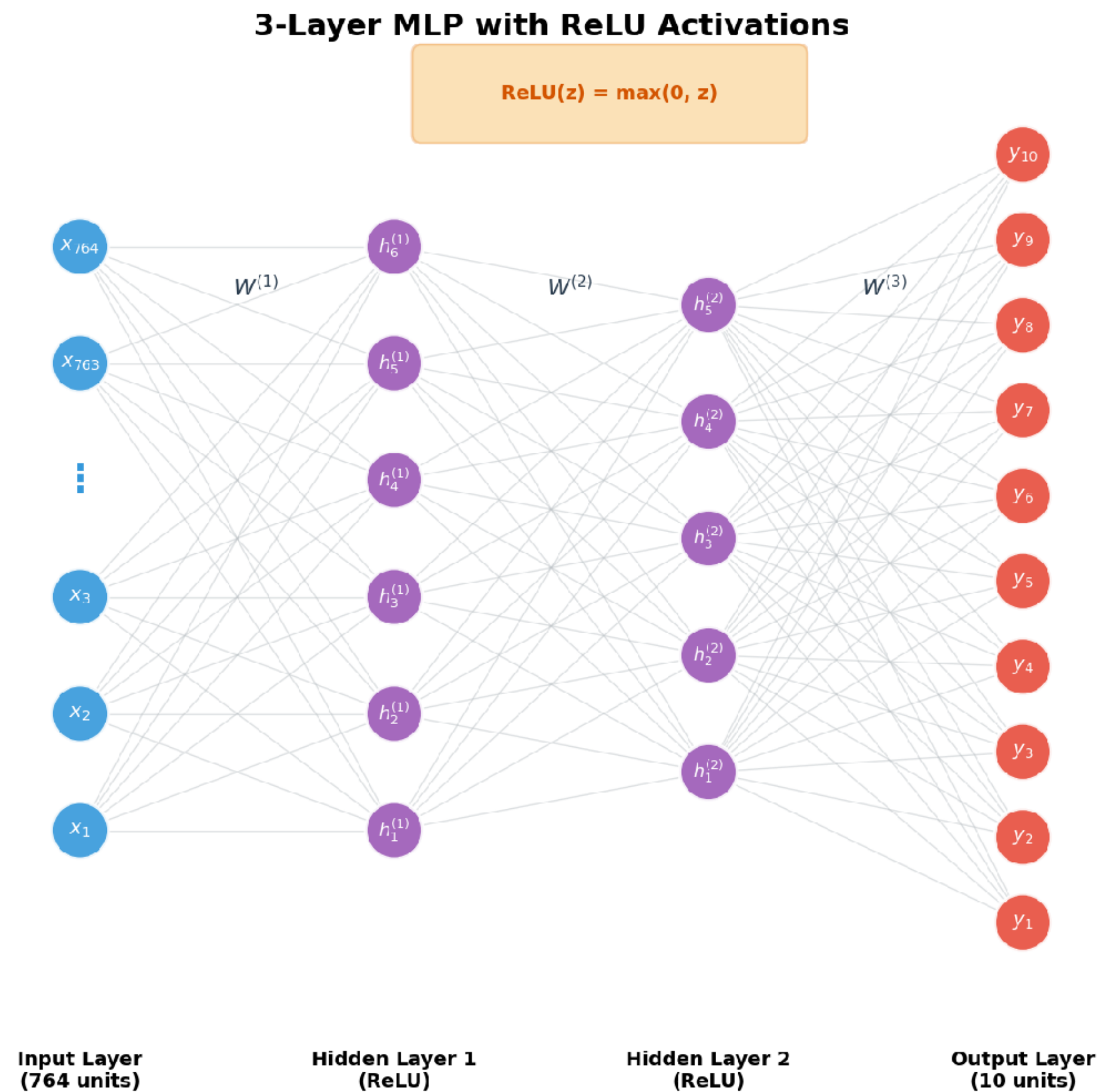
# Your First Deep Network: Why?



3-Layer MLP with ReLU Activations

$ReLU(z) = max(0, z)$

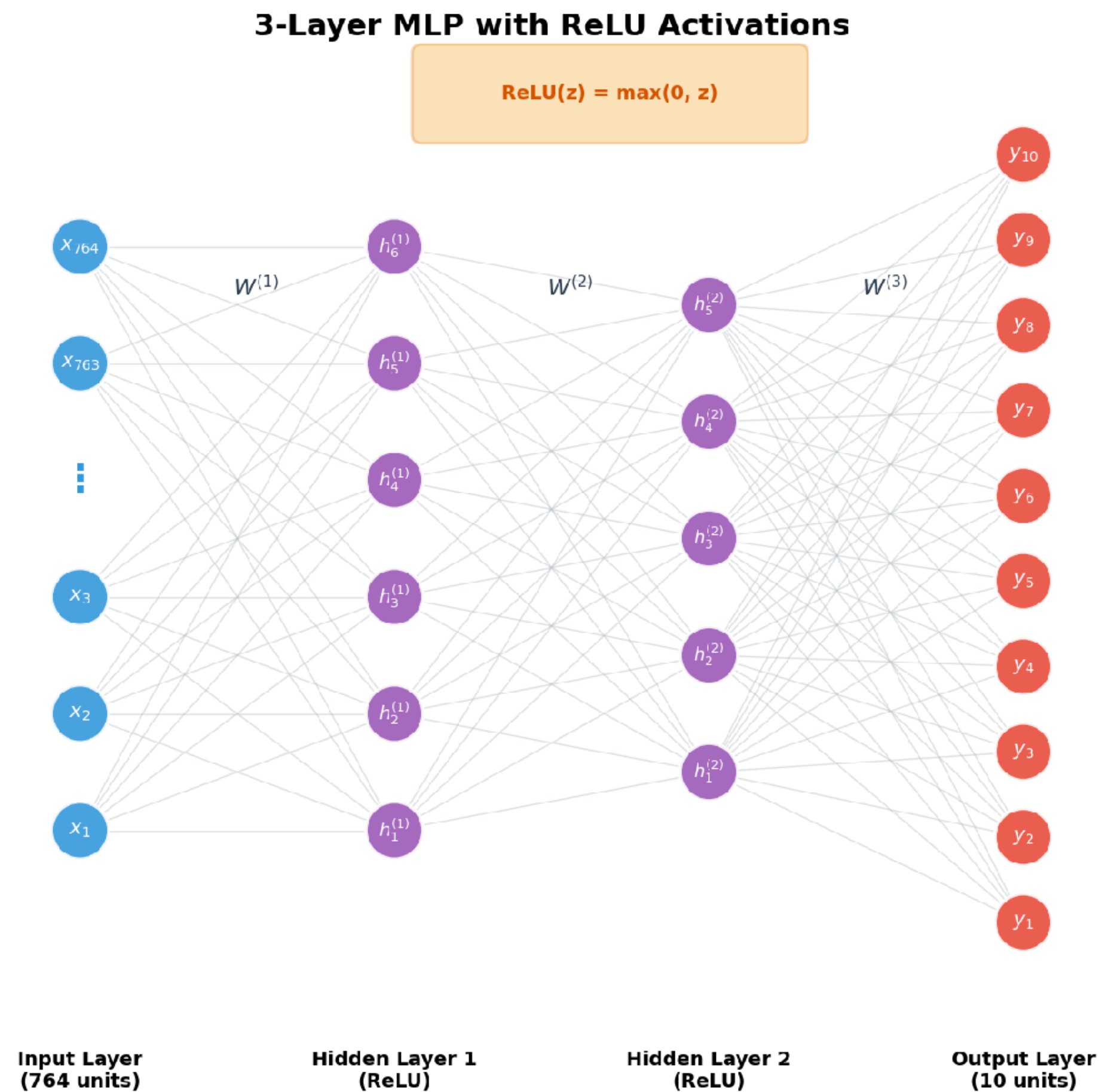Input Layer (764 units) · Hidden Layer 1 (ReLU) · Hidden Layer 2 (ReLU) · Output Layer (10 units)

- Why 2 hidden layers?

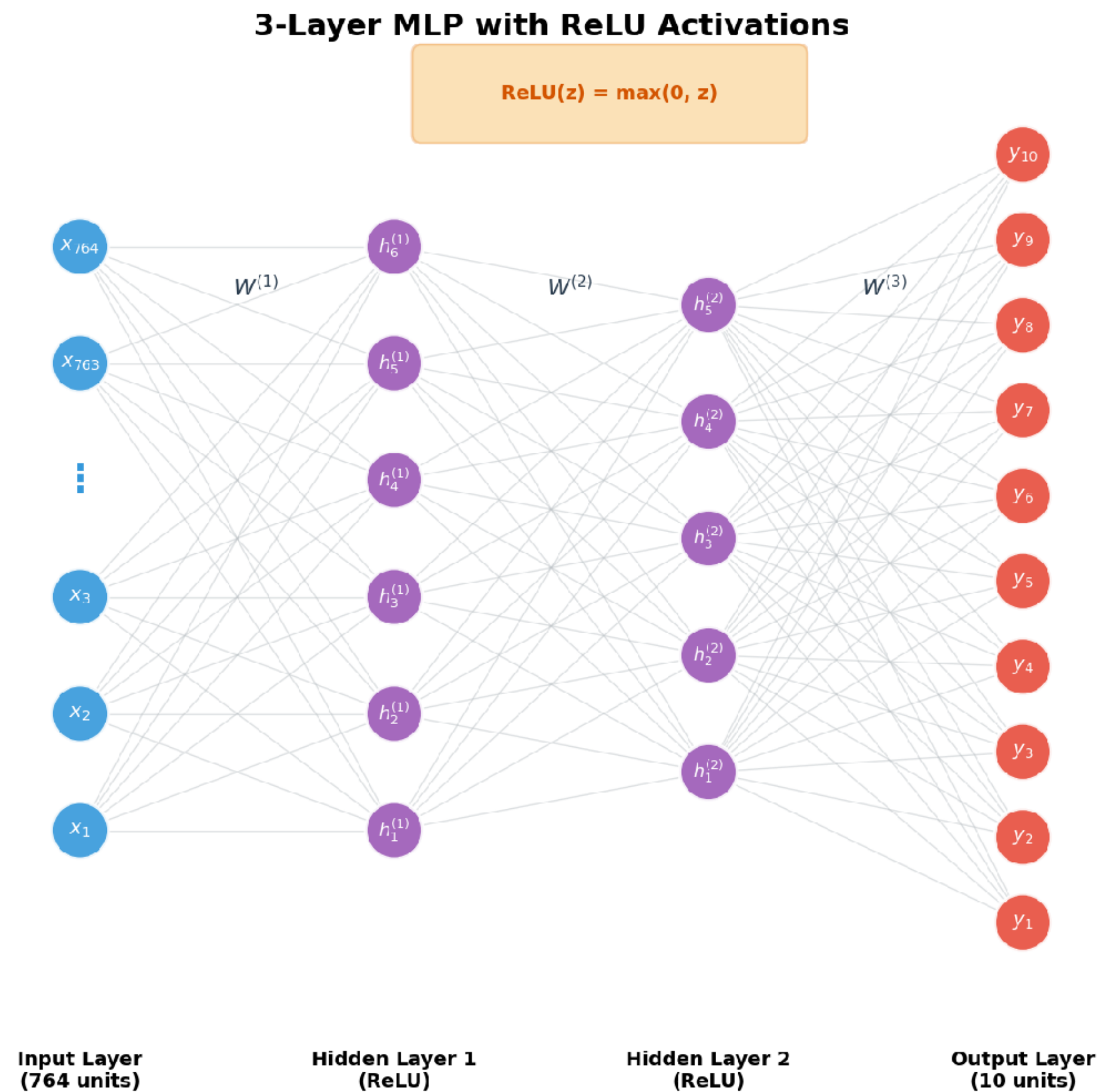- Why do we need a "ReLU"?

# Your First Deep Network: Why?



**3-Layer MLP with ReLU Activations**

$$\text{ReLU}(z) = \max(0, z)$$

$x_{764}$   $x_{763}$   $x_3$   $x_2$   $x_1$

$h_6^{(1)}$   $h_5^{(1)}$   $h_4^{(1)}$   $h_3^{(1)}$   $h_2^{(1)}$   $h_1^{(1)}$

$h_5^{(2)}$   $h_4^{(2)}$   $h_3^{(2)}$   $h_2^{(2)}$   $h_1^{(2)}$

$y_{10}$   $y_9$   $y_8$   $y_7$   $y_6$   $y_5$   $y_4$   $y_3$   $y_2$   $y_1$

$W^{(1)}$   $W^{(2)}$   $W^{(3)}$

**Input Layer (764 units)**   **Hidden Layer 1 (ReLU)**   **Hidden Layer 2 (ReLU)**   **Output Layer (10 units)**
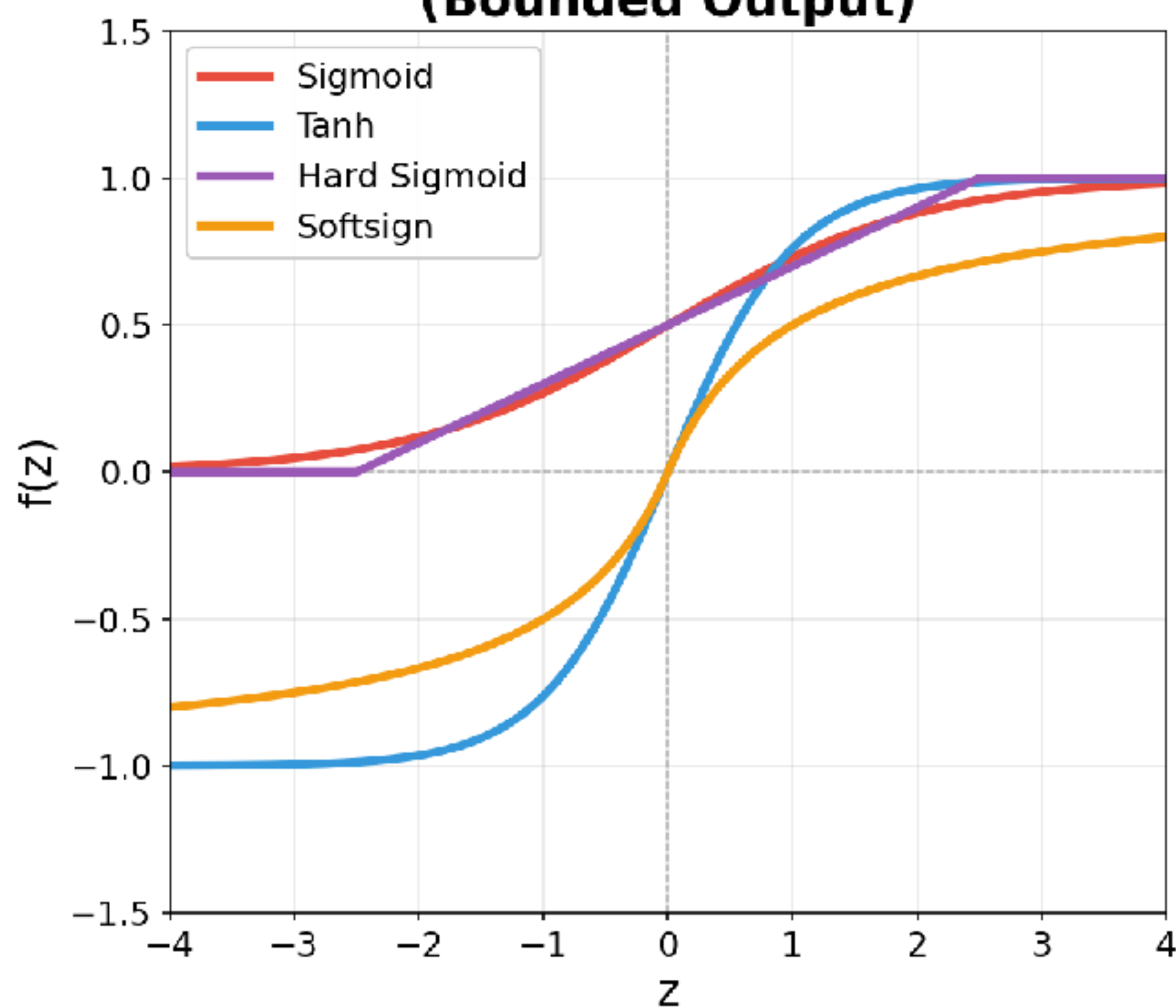
- Why 2 hidden layers?

- Why do we need a "ReLU"?

- What else besides a "ReLU"?
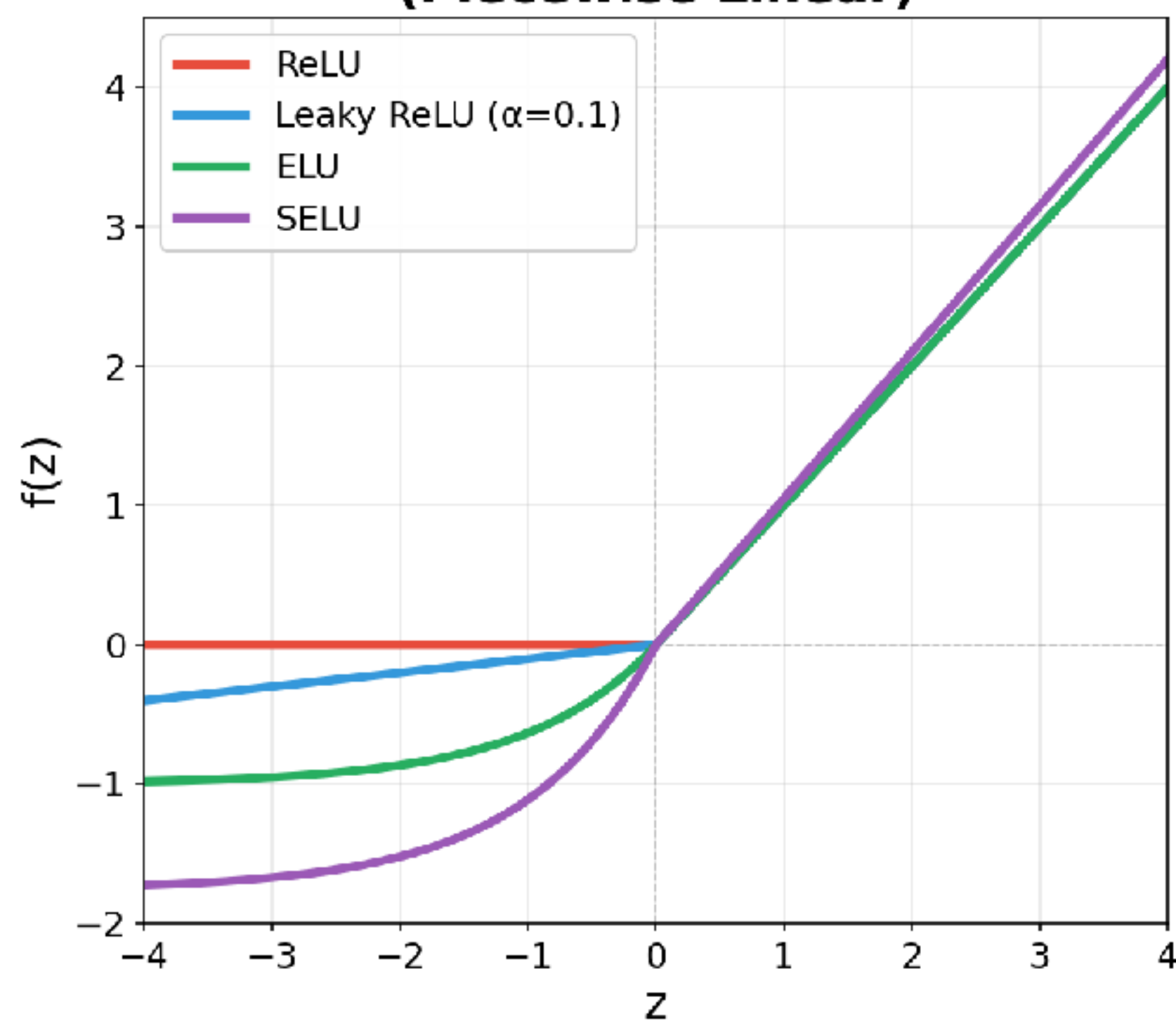
# Your First Deep Network: Why?



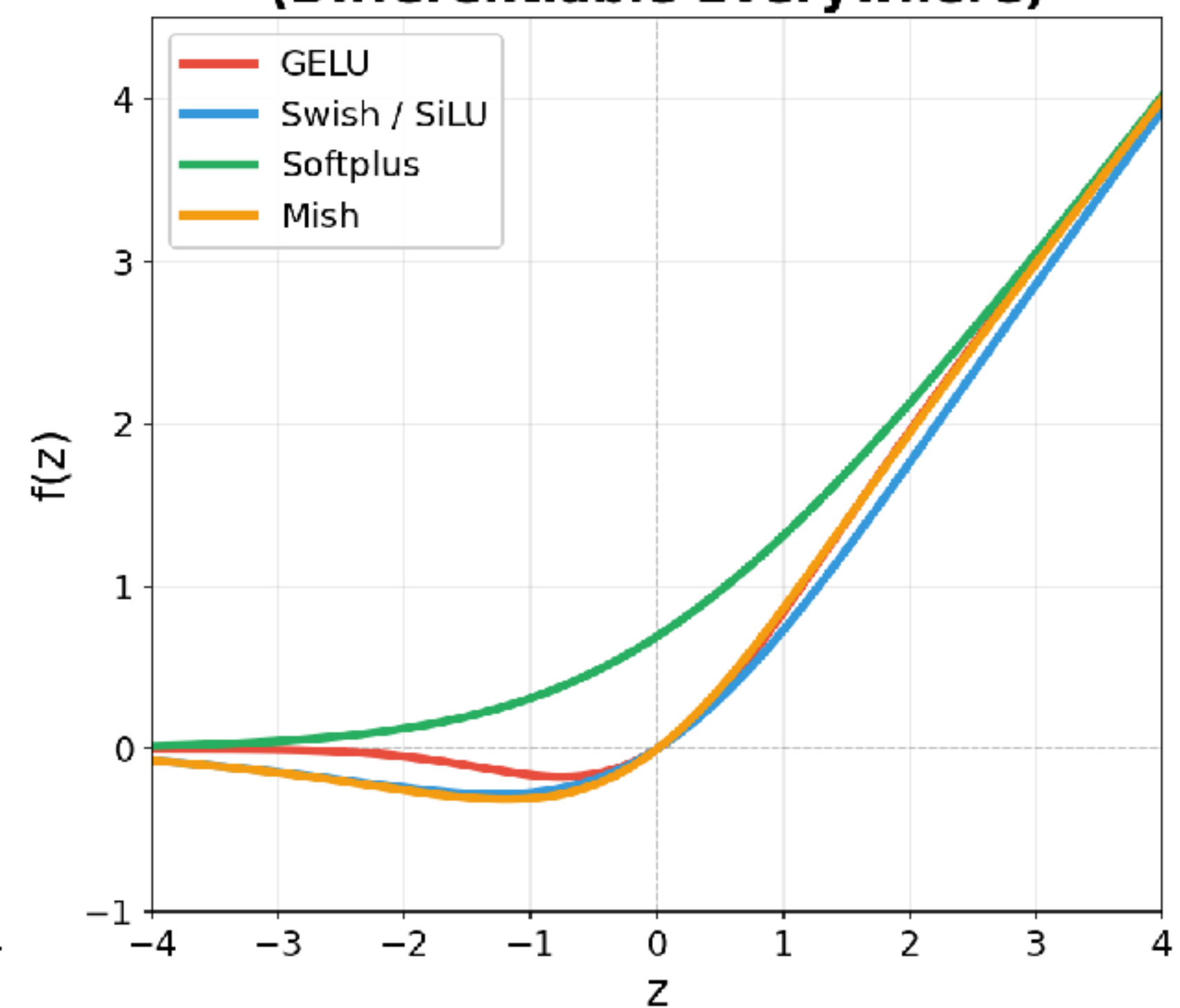Common Activation Functions in Neural Networks

# Your First Deep Network: Why?



Common Activation Functions in Neural Networks
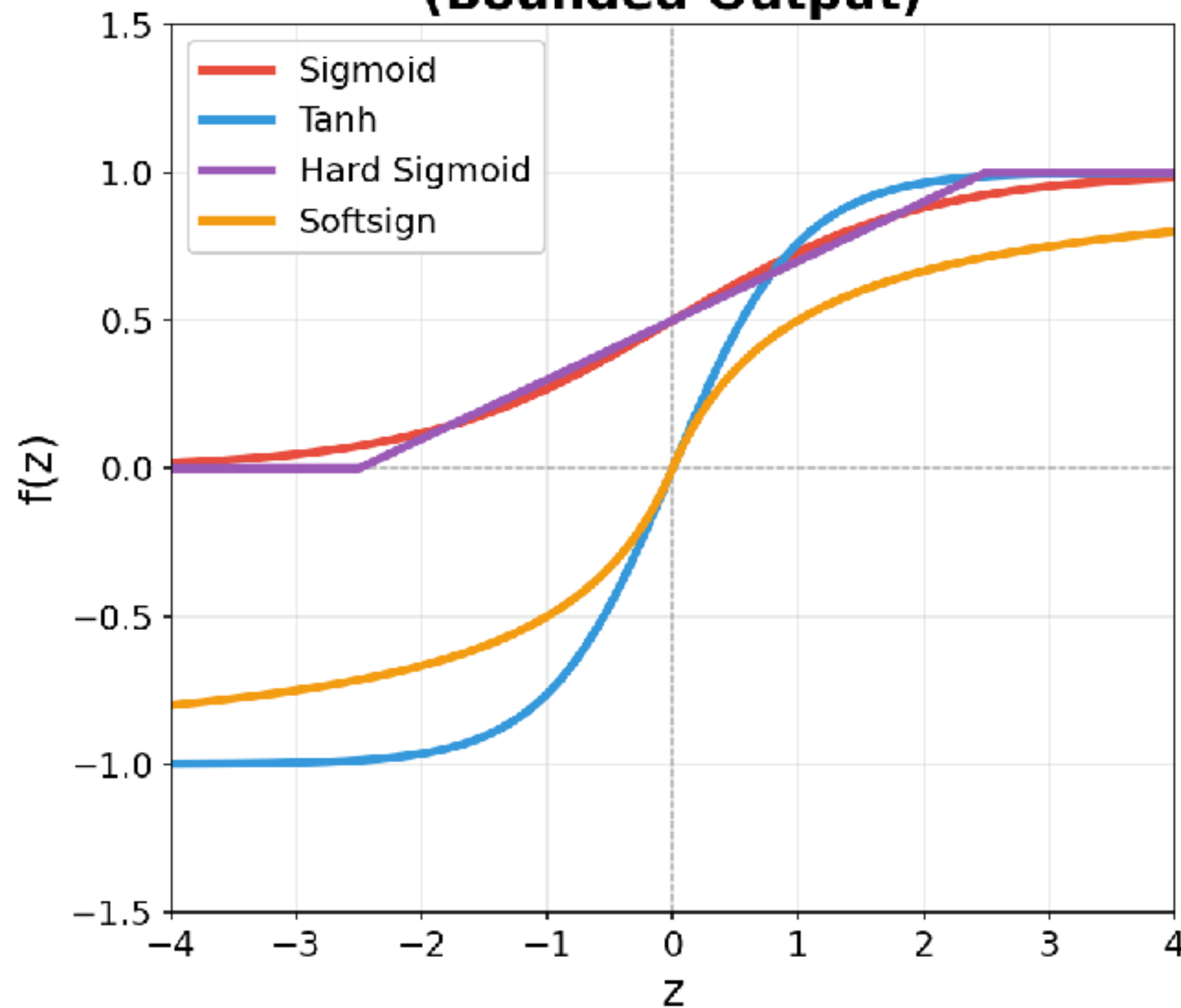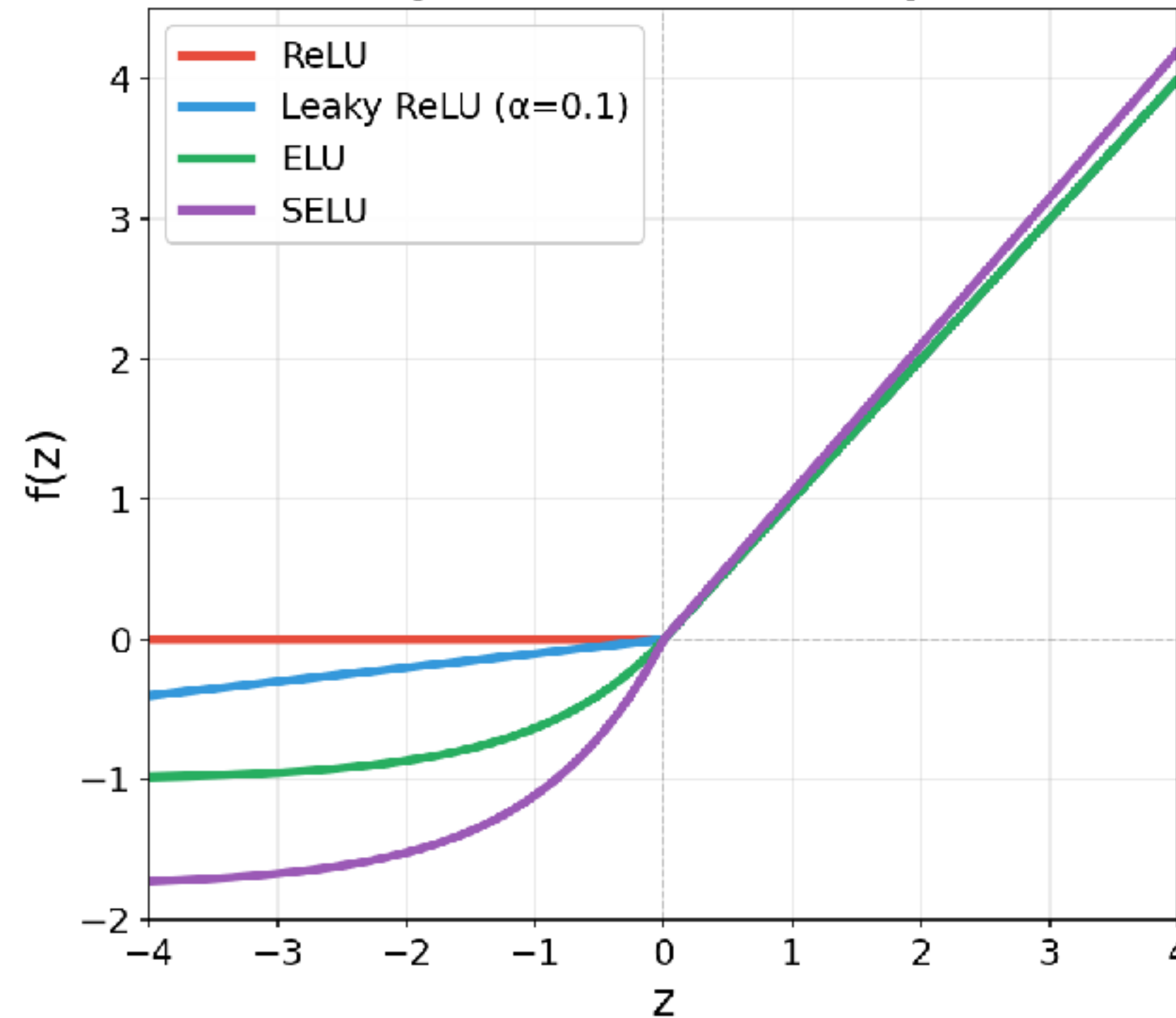
**Sigmoid-like (Bounded Output)** — Sigmoid, Tanh, Hard Sigmoid, Softsign

**ReLU Family (Piecewise Linear)** — ReLU, Leaky ReLU ($\alpha=0.1$), ELU, SELU

**Smooth Modern (Differentiable Everywhere)** — GELU, Swish / SiLU, Softplus, Mish

Can you think of "problems" for some of them?

# Your First Deep Network!



3-Layer MLP with ReLU Activations

$$ReLU(z) = max(0, z)$$

Input Layer
(764 units)

Hidden Layer 1
(ReLU)

Hidden Layer 2
(ReLU)

Output Layer
(10 units)

# Your First Deep Network!



## 3-Layer MLP with ReLU Activations

ReLU(z) = max(0, z)

$W^{(1)}$  $W^{(2)}$  $W^{(3)}$

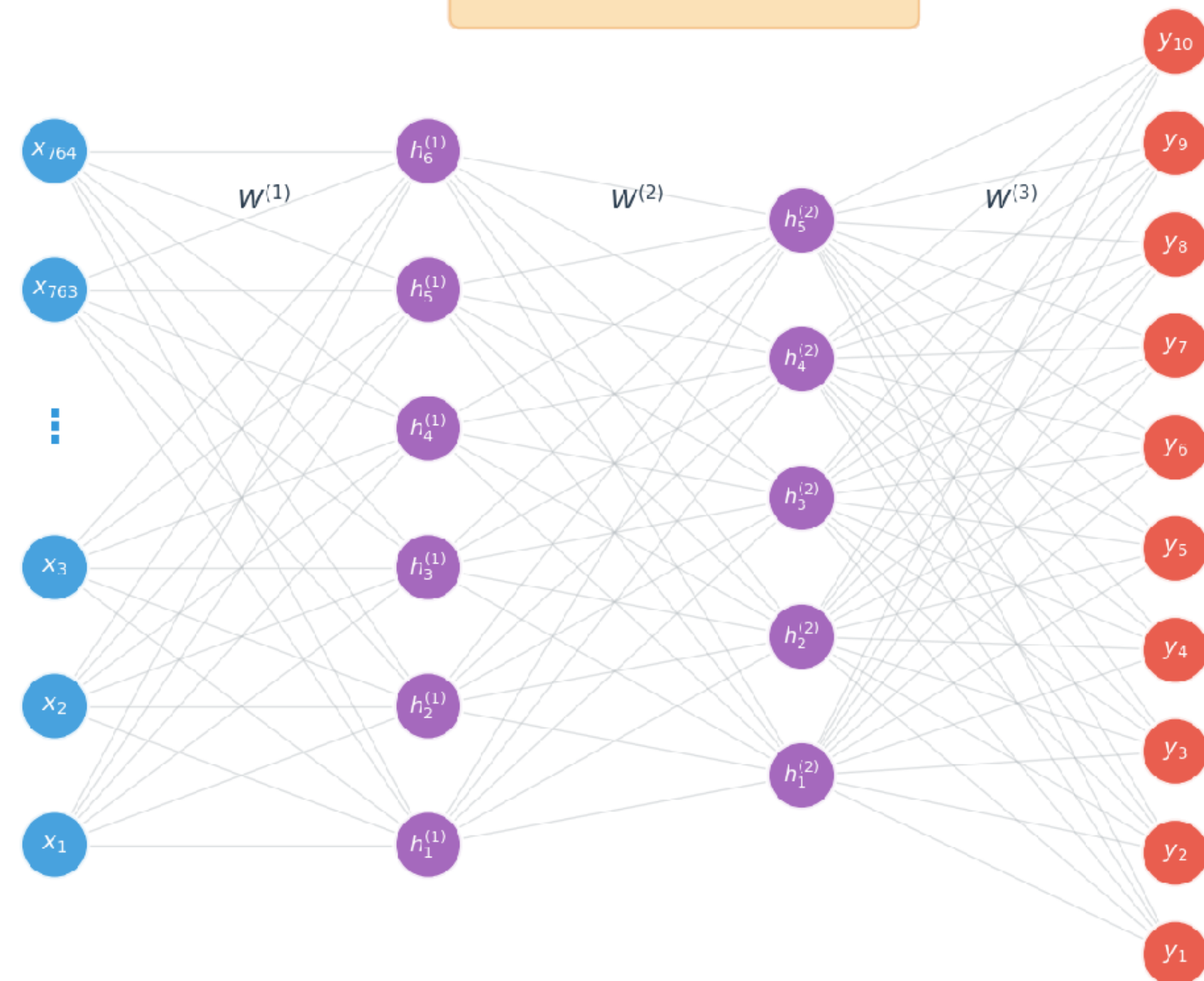$x_{764}$  $x_{763}$  $x_3$  $x_2$  $x_1$

$h_6^{(1)}$  $h_5^{(1)}$  $h_4^{(1)}$  $h_3^{(1)}$  $h_2^{(1)}$  $h_1^{(1)}$

$h_5^{(2)}$  $h_4^{(2)}$  $h_3^{(2)}$  $h_2^{(2)}$  $h_1^{(2)}$

$y_{10}$  $y_9$  $y_8$  $y_7$  $y_6$  $y_5$  $y_4$  $y_3$  $y_2$  $y_1$

**Input Layer
(764 units)**  **Hidden Layer 1
(ReLU)**  **Hidden Layer 2
(ReLU)**  **Output Layer
(10 units)**

What training loss?

# Your First Deep Network!



**3-Layer MLP with ReLU Activations**

$ReLU(z) = max(0, z)$

What training loss?

$p(y = 9 \,|\, x_n)$

Input Layer
(764 units)

Hidden Layer 1
(ReLU)

Hidden Layer 2
(ReLU)

Output Layer
(10 units)

# Your First Deep Network!

**3-Layer MLP with ReLU Activations**

ReLU(z) = max(0, z)



Input Layer
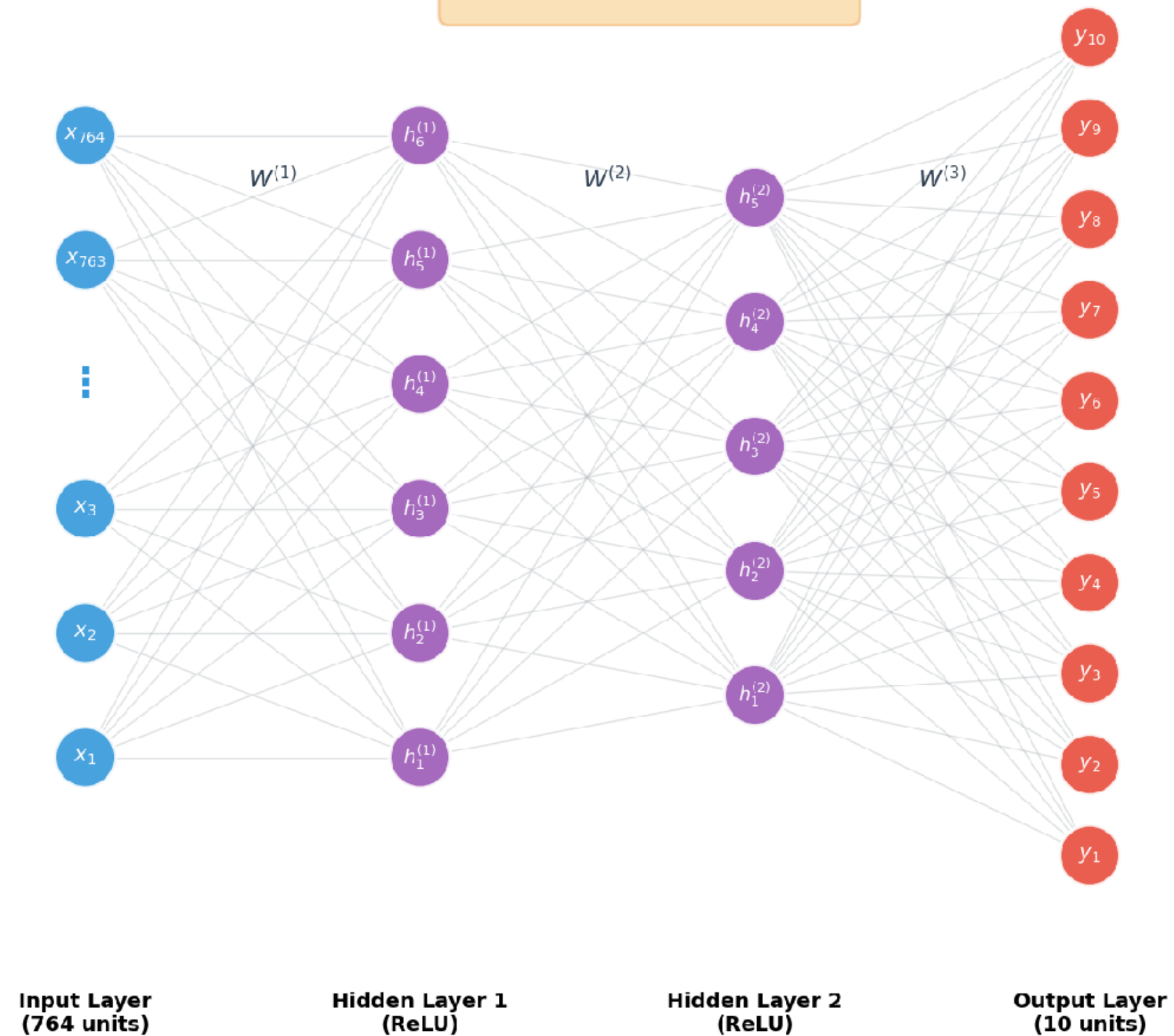(764 units)
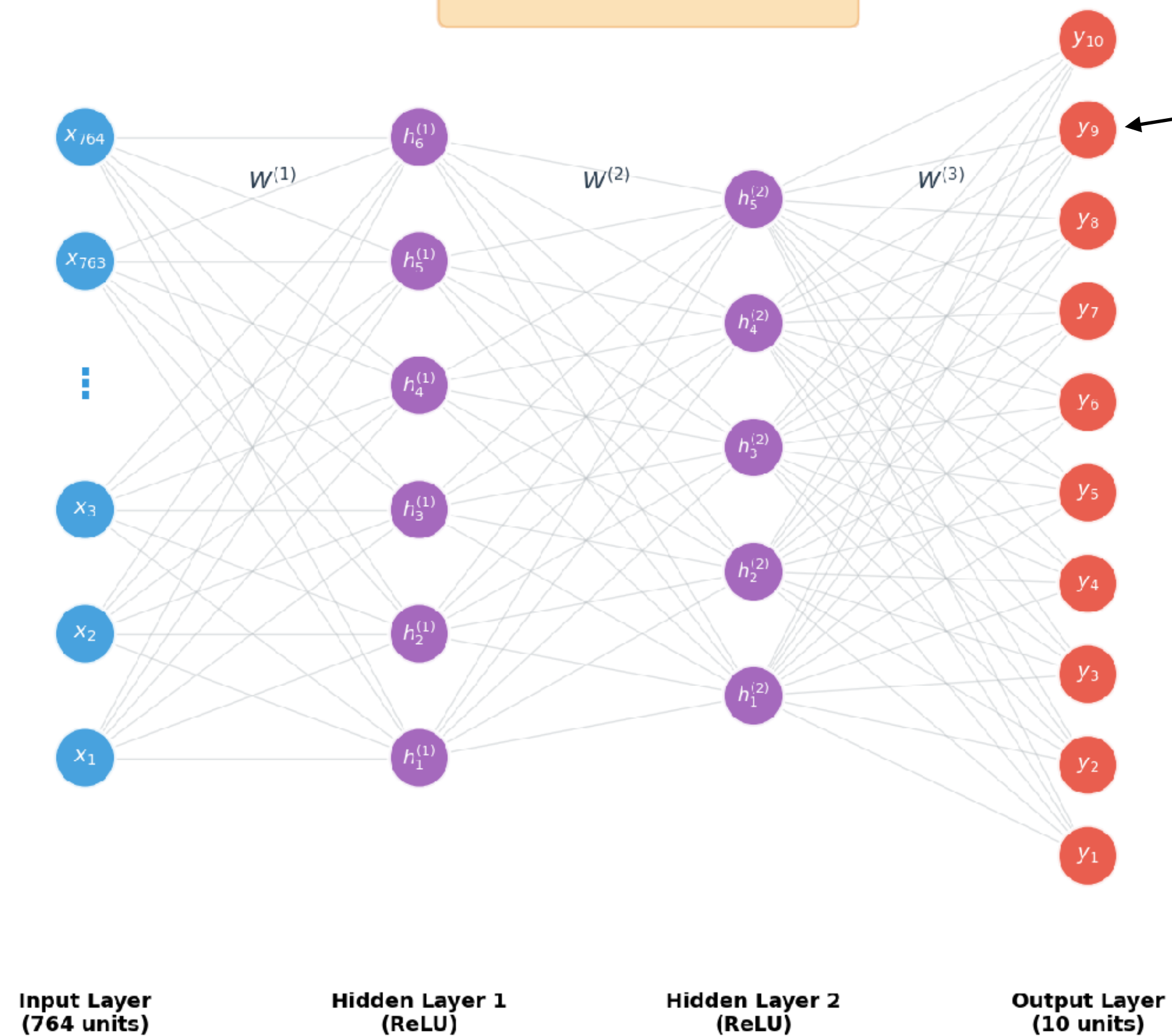
Hidden Layer 1
(ReLU)

Hidden Layer 2
(ReLU)

Output Layer
(10 units)

**What training loss?**

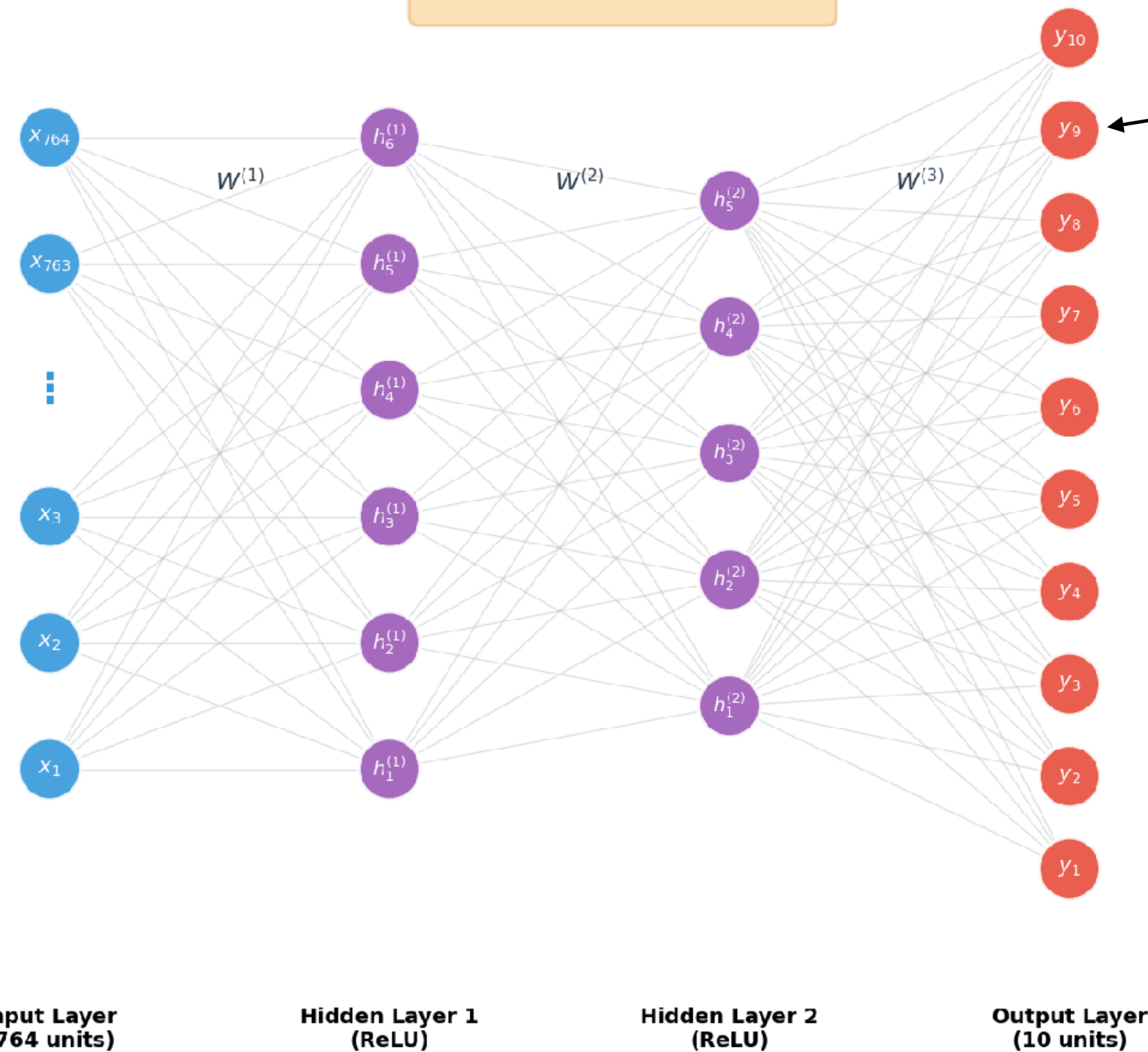$$p(y = 9 \mid x_n)$$

$$\mathscr{L} = -\sum_{n=1}^{N} \sum_{c=1}^{10} 1_{\{y_n = c\}} \log(p(y = c \mid x_n))$$

# Your First Deep Network!

**3-Layer MLP with ReLU Activations**

$$\text{ReLU}(z) = \max(0, z)$$



**Input Layer**
**(764 units)**

**Hidden Layer 1**
**(ReLU)**

**Hidden Layer 2**
**(ReLU)**

**Output Layer**
**(10 units)**

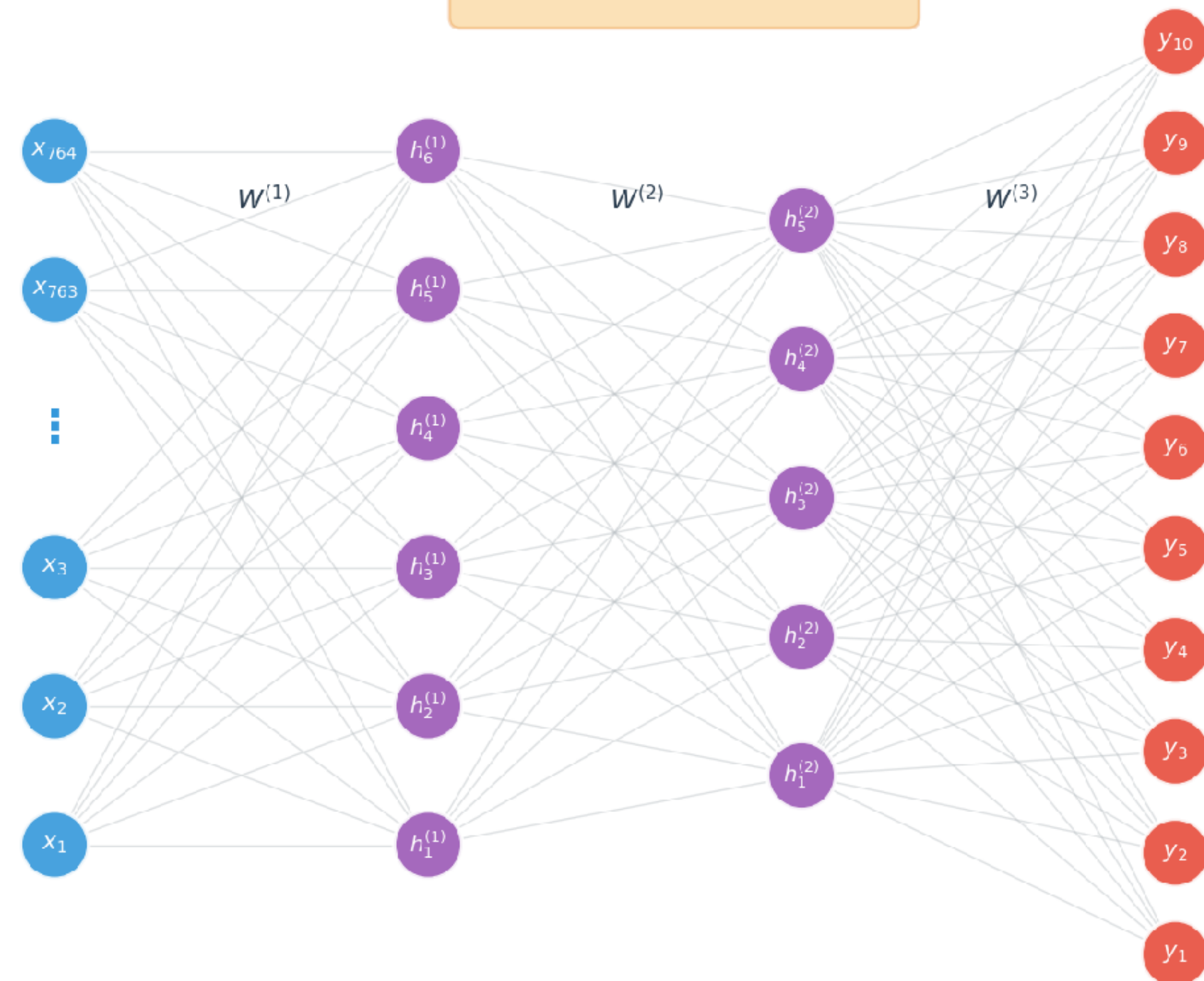What training loss?

$$p(y = 9 \mid x_n)$$

$$\mathcal{L} = -\sum_{n=1}^{N} \sum_{c=1}^{10} 1_{\{y_n = c\}} \log(p(y = c \mid x_n))$$

$$\hat{y}_n = \arg \max_{c=1,\dots,10} p(y = c \mid x_n)$$

# Your First Deep Network: Training



3-Layer MLP with ReLU Activations

ReLU(z) = max(0, z)

$x_{764}$     $x_{763}$     $x_3$     $x_2$     $x_1$

$W^{(1)}$     $W^{(2)}$     $W^{(3)}$

$h_6^{(1)}$, $h_5^{(1)}$, $h_4^{(1)}$, $h_3^{(1)}$, $h_2^{(1)}$, $h_1^{(1)}$

$h_5^{(2)}$, $h_4^{(2)}$, $h_3^{(2)}$, $h_2^{(2)}$, $h_1^{(2)}$

$y_{10}$, $y_9$, $y_8$, $y_7$, $y_6$, $y_5$, $y_4$, $y_3$, $y_2$, $y_1$
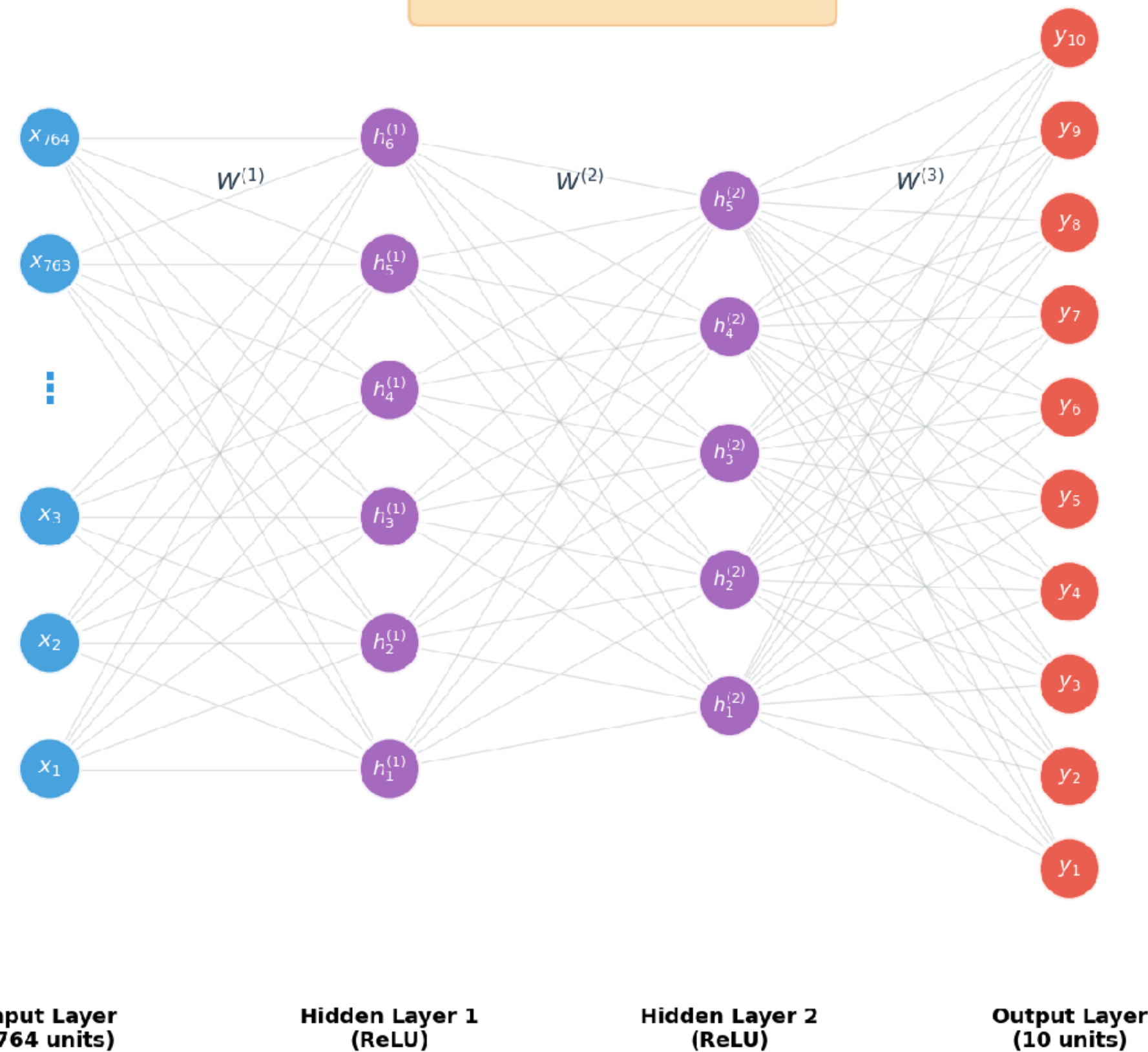
Input Layer
(764 units)

Hidden Layer 1
(ReLU)

Hidden Layer 2
(ReLU)

Output Layer
(10 units)

# Your First Deep Network: Training



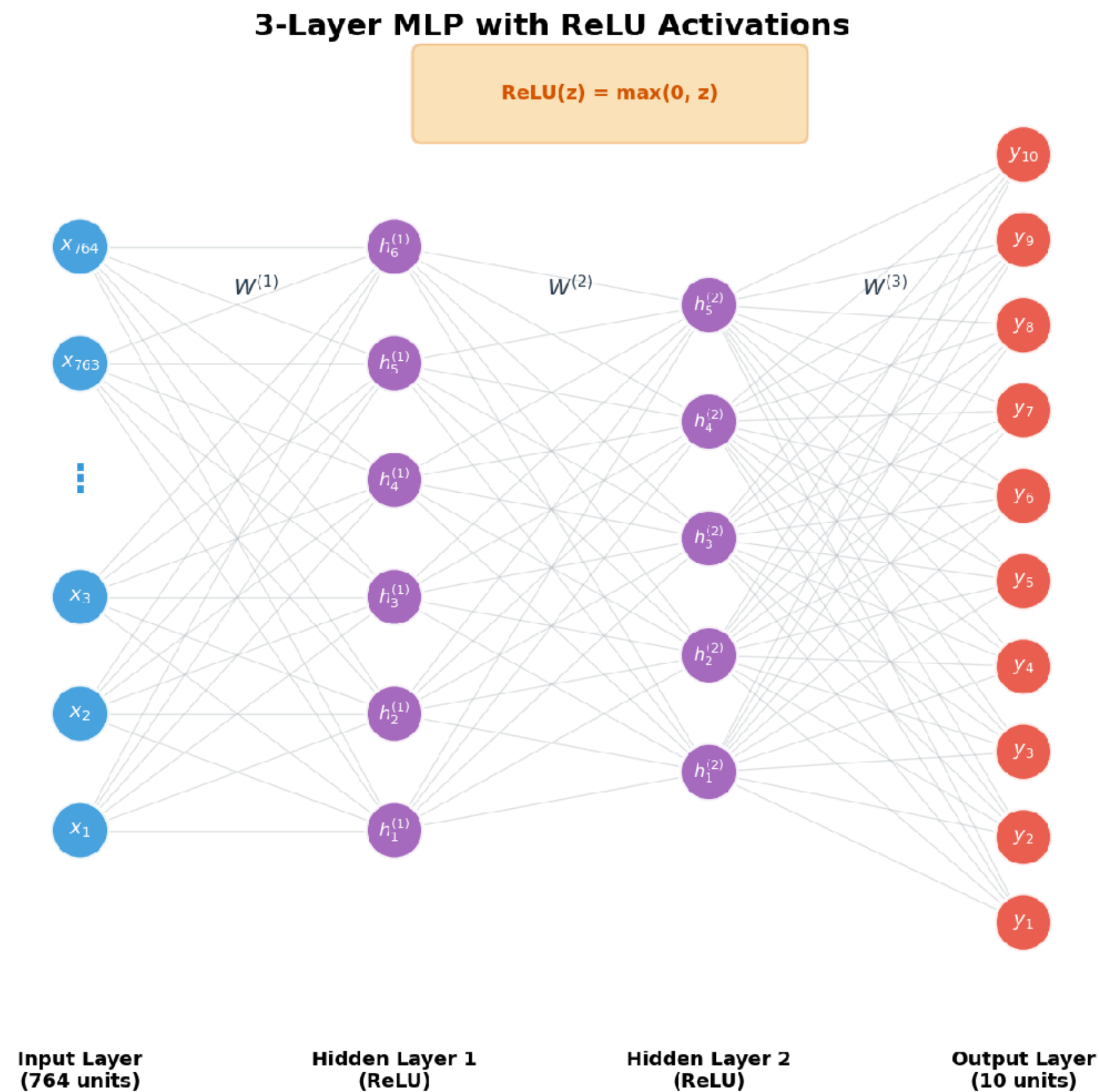**3-Layer MLP with ReLU Activations**

ReLU(z) = max(0, z)

$x_{764}$ $x_{763}$ $x_3$ $x_2$ $x_1$

$h_6^{(1)}$ $h_5^{(1)}$ $h_4^{(1)}$ $h_3^{(1)}$ $h_2^{(1)}$ $h_1^{(1)}$

$W^{(1)}$

$h_5^{(2)}$ $h_4^{(2)}$ $h_3^{(2)}$ $h_2^{(2)}$ $h_1^{(2)}$

$W^{(2)}$

$W^{(3)}$

$y_{10}$ $y_9$ $y_8$ $y_7$ $y_6$ $y_5$ $y_4$ $y_3$ $y_2$ $y_1$

**Input Layer (764 units)**  **Hidden Layer 1 (ReLU)**  **Hidden Layer 2 (ReLU)**  **Output Layer (10 units)**

How to train the parameters $W^{(1)}, W^{(2)}, W^{(3)}$?

# Your First Deep Network: Training



### 3-Layer MLP with ReLU Activations

$ReLU(z) = max(0, z)$

Input Layer (764 units) — Hidden Layer 1 (ReLU) — Hidden Layer 2 (ReLU) — Output Layer (10 units)
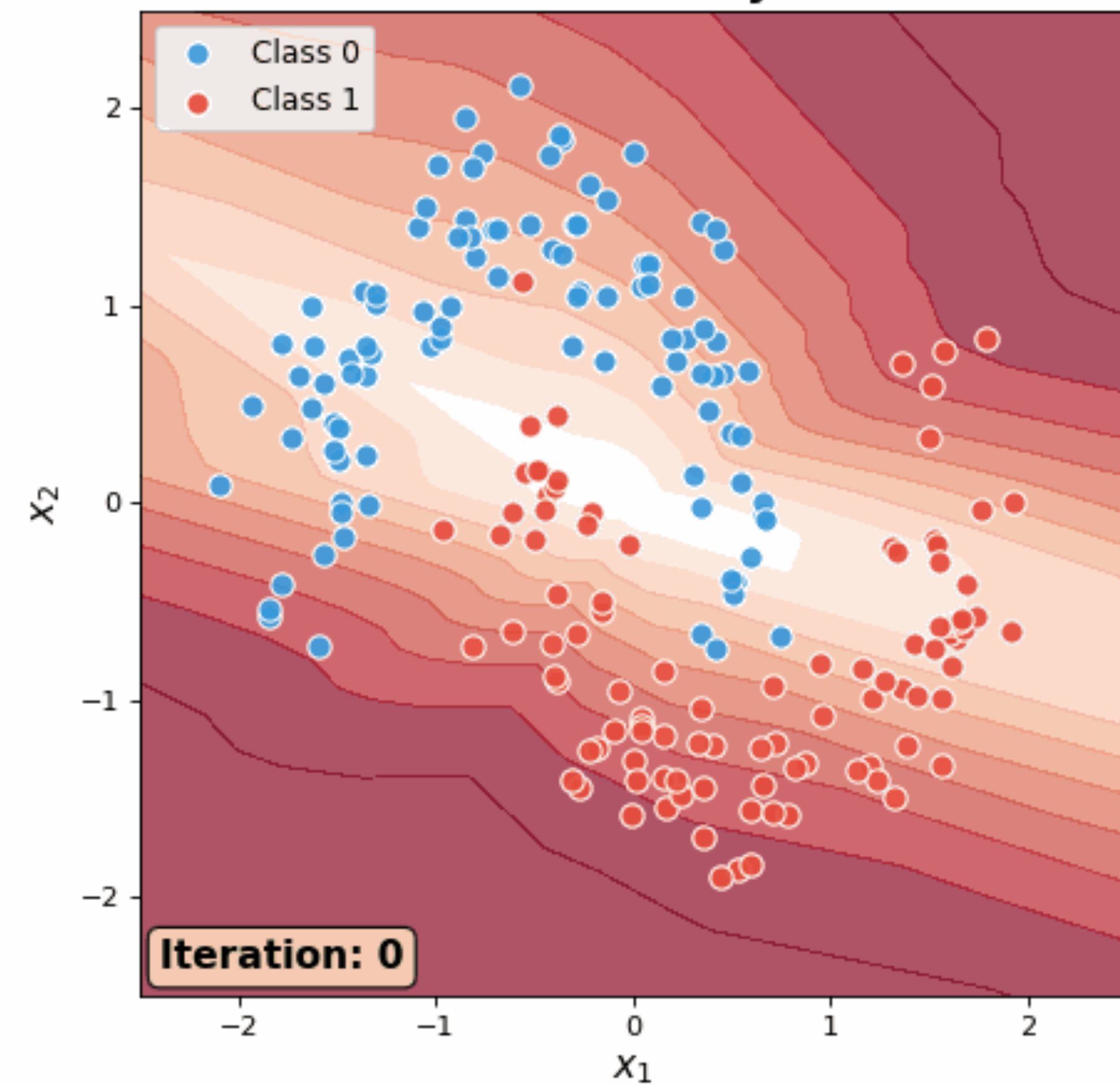
How to train the parameters $W^{(1)}, W^{(2)}, W^{(3)}$?

*Same old gradient descent!*

*At home: try to derive the gradient for those 3 matrices*

# Your First Deep Network: Action!

# Your First Deep Network: Action!
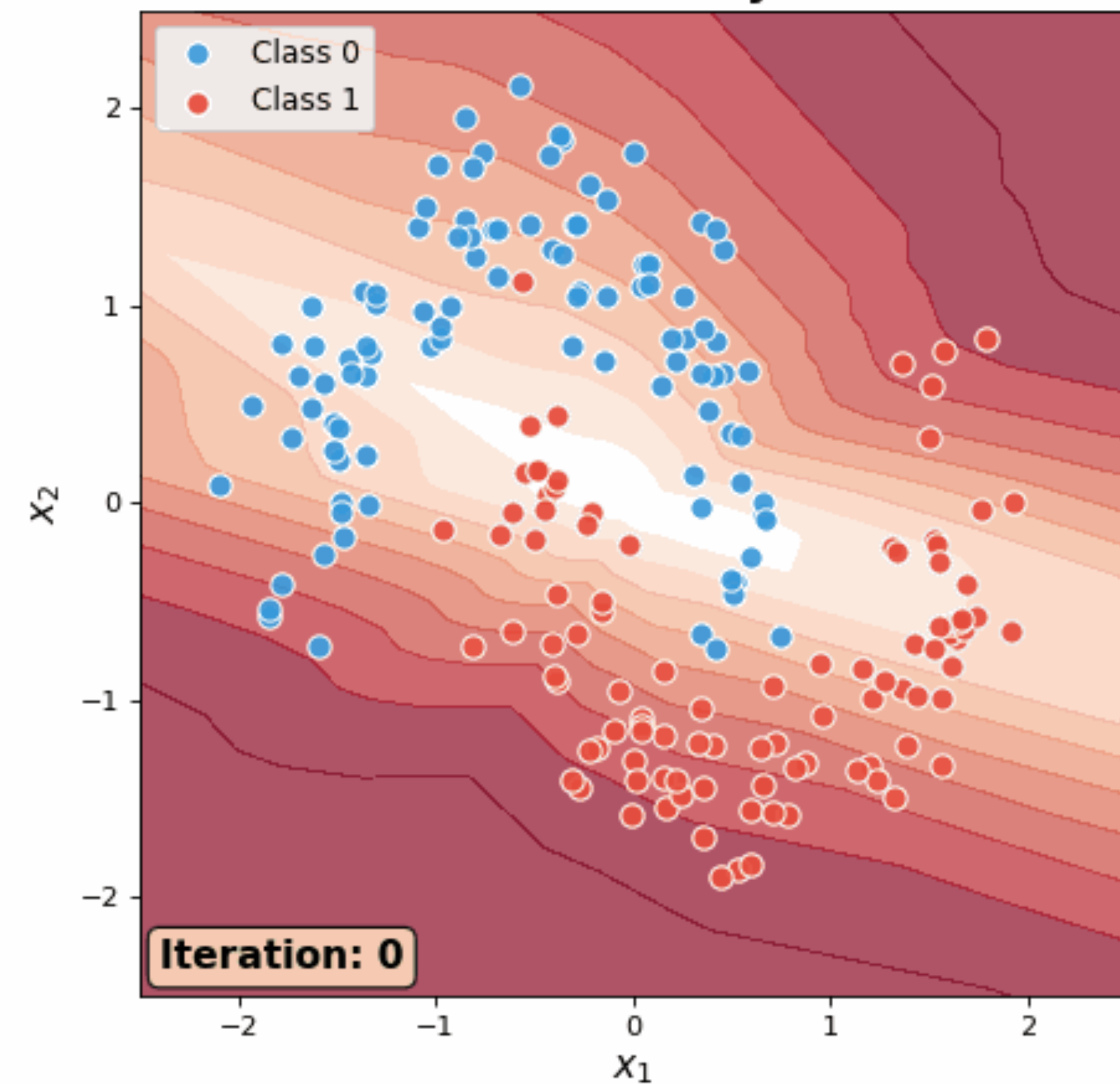
# Open Discussion

# Open Discussion

- Can a Deep Network solve anything?

# Open Discussion

- Can a Deep Network solve anything?

- How to search for the right "architecture"?

# Open Discussion

- Can a Deep Network solve anything?

- How to search for the right "architecture"?

- How many lines of codes to implement the MNIST model and reach 99.5%?

# Questions?