

**APL 745**

# Supervised Learning: Feedforward neural network

*Dr. Rajdip Nayek*

Block V, Room 418D

Department of Applied Mechanics

Indian Institute of Technology Delhi

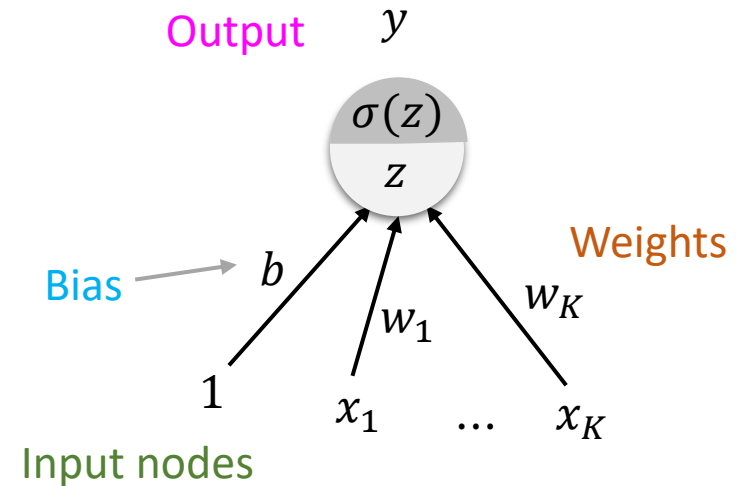
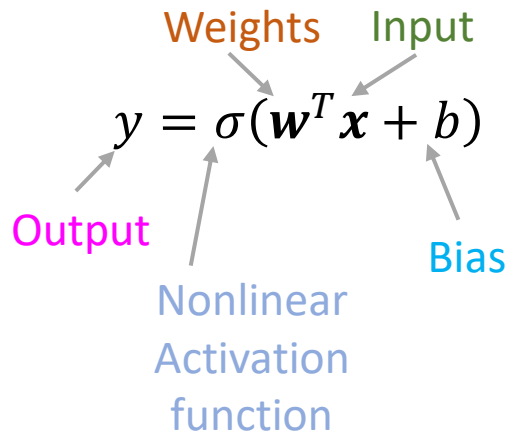
E-mail: [rajdipn@am.iitd.ac.in](mailto:rajdipn@am.iitd.ac.in)

# Learning goals

- Know the basic terminology for neural networks
- Given the weights and biases for a neural network, be able to compute its output from its input
- Know the types of activation functions used
- Be able to hand-design the weights of a neural net to represent functions like XOR
- Understand why shallow neural nets are expressive, and why this isn't necessarily very interesting

# Overview

- In the first lecture, we introduced a neuron-like processing unit



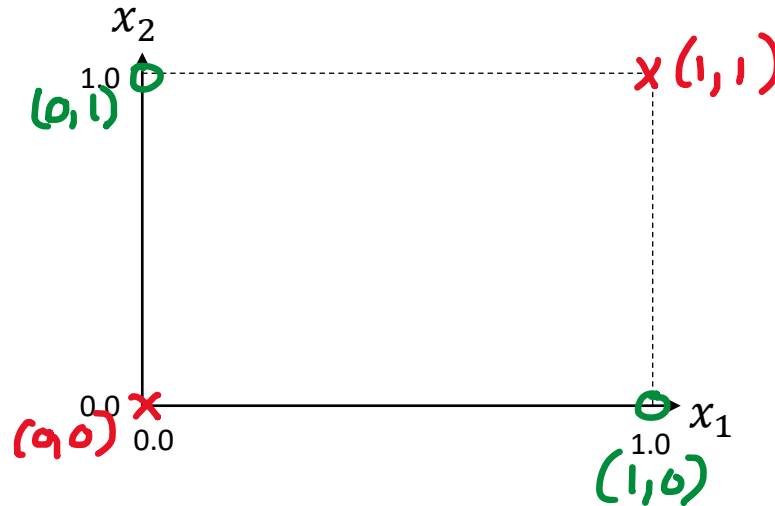
➤ For **linear regression**:  $\sigma(z) = z$

➤ For **linear classification**:  $\sigma(z) = \left\{ \begin{array}{l} \text{Step function} \rightarrow \text{Perceptron} \\ \text{Sigmoid curve} \rightarrow \text{Sigmoid neuron} \\ \text{Softmax curve} \rightarrow \text{Softmax for multi-class} \end{array} \right.$

- A **neural network** is just a combination of lots of these units

# Overview

- We noted that **there are functions that can't be represented by linear models**
  - Ex: Linear regression can not represent quadratic functions
  - Ex: Linear classifiers can not represent the XOR gate



- We also saw one particular way of avoiding the above issue
  - By defining **nonlinear features** or basis functions  $\phi(x)$

# Limitations of linear models

- While linear models are convenient, the way of defining nonlinear features is not very satisfying
  - The nonlinear features  $\phi(x)$  have to be chosen and fixed manually and that requires expert domain knowledge
  - It might require a large number of nonlinear features to represent certain functions
- **Question:** Can we work with unrestricted non-linear models?
  - We don't want to work with a finite set of basis functions
  - We don't want to completely specify the basis functions beforehand
- **Answer:** Yes, now we are going to see how we can free ourselves from this restriction

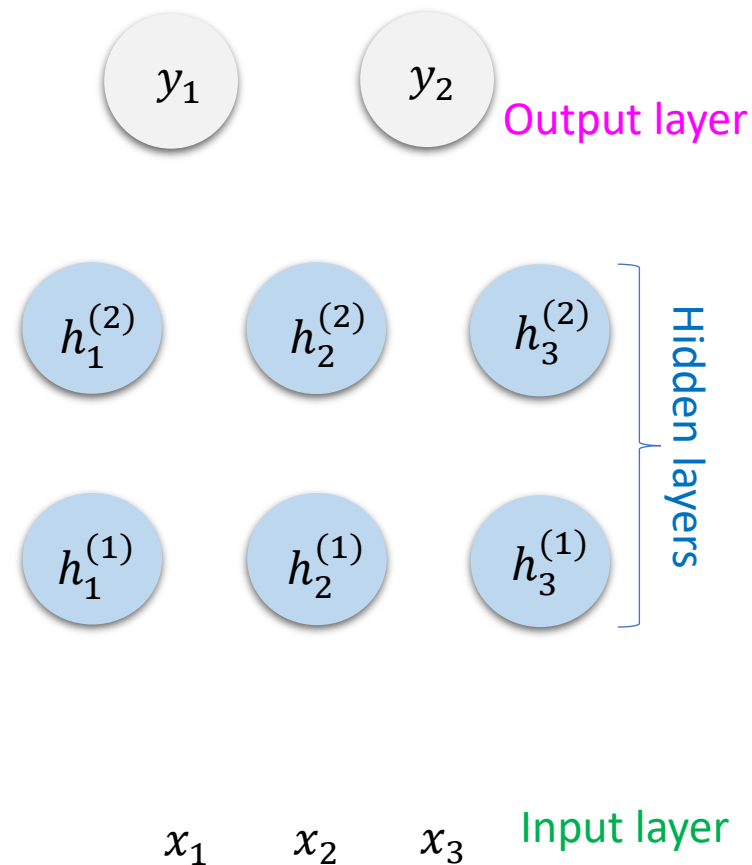
# Flexible nonlinear models

- **Idea 1:** Instead of choosing a fixed set of basis functions apriori, we can choose our **basis functions to depend on the data** and then allow a very large number (sometimes an infinite number) of them to be included without paying a computational price. This kind approach falls under **kernel** methods (e.g. **Support Vector Machines**)
  - Was a popular idea in the 2000s upto around 2010
- **Idea2:** Deep learning → **Learn non-linear basis functions** (e.g. **Multilayer Neural Networks**)
  - Became very popular from around 2010
  - Even though the number of basis functions are finite and fixed, the parameters of the basis functions are adapted (or learned) so as to effectively vary them based on the data

# A multilayer neural network

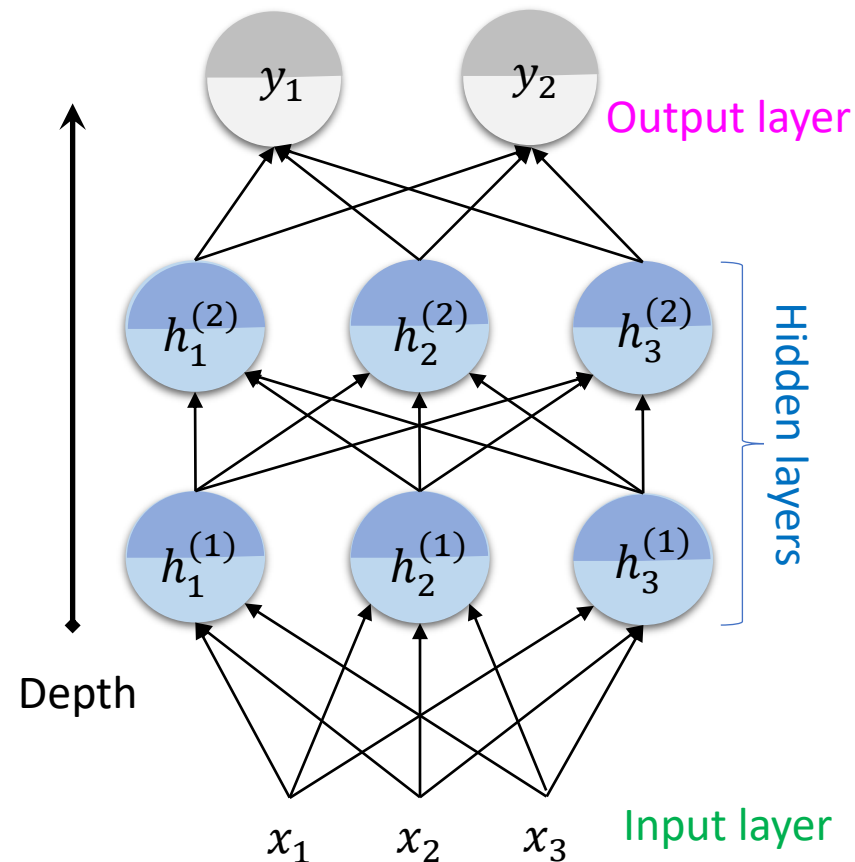
## Terminology

- **Multilayer neural network** is a network of neuron-like units arranged in **layers**. This network in picture has 4 layers
- The layer containing the input features  $(x_1, x_2, x_3)$  is called the **input layer (Layer 0)**
- The two middle layers containing 3 units are called **hidden layers**
  - The outputs of the first hidden layer are denoted by  $h_1^{(1)}, h_2^{(1)}, h_3^{(1)}$
  - The outputs of the second hidden layer are denoted by  $h_1^{(2)}, h_2^{(2)}, h_3^{(2)}$
- The final layer is called the **output layer** (here two output units are shown)



# A multilayer neural network

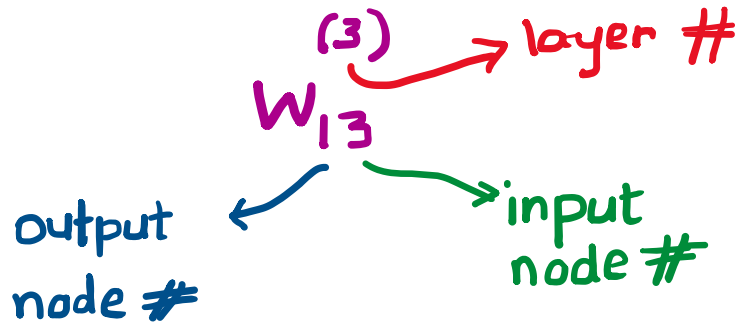
- Each neuron in the hidden layers and the output layer can be split into two parts: **pre-activation** and **activation**
- Fully connected** ← every unit in one layer is connected to every unit in the next layer
- Depth** ← the number of computational layers
  - Three-layered network
  - Input layer is layer 0
- Width** ← the number of units in a layer
  - Three-unit wide hidden layers (each one could have different width as well)
  - Three-unit wide input layer
  - Two-unit wide output layer
- Information flows in only forward direction, from input nodes through hidden nodes to output nodes, hence called **feed-forward neural network (FNN)**. No cycles or loops



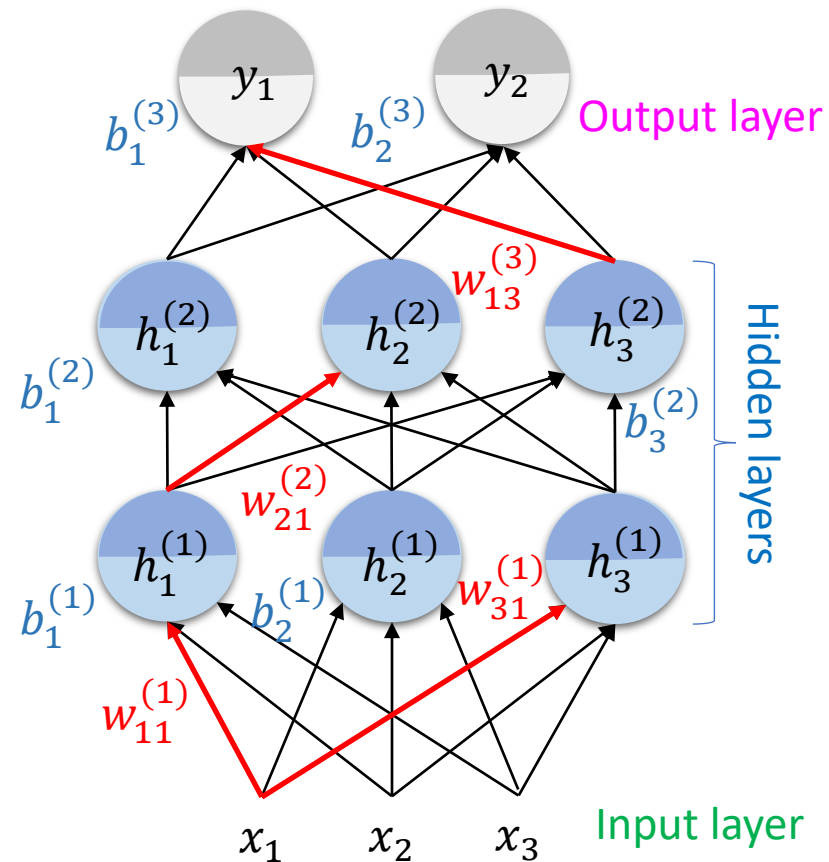
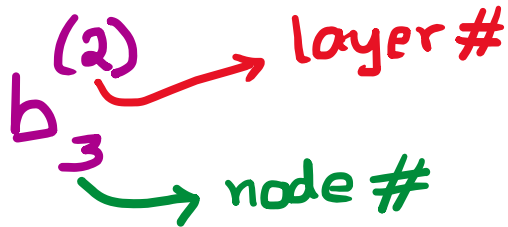


# Computation in an FNN

- Each link (arrow) is associated with a weight
  - $w_{ij}^{(k)} \leftarrow k^{\text{th}} \text{ layer}, i^{\text{th}} \text{ output}, j^{\text{th}} \text{ input}$



- Each circular node has its own bias
  - $b_i^{(k)} \leftarrow k^{\text{th}} \text{ layer}, i^{\text{th}} \text{ unit}$



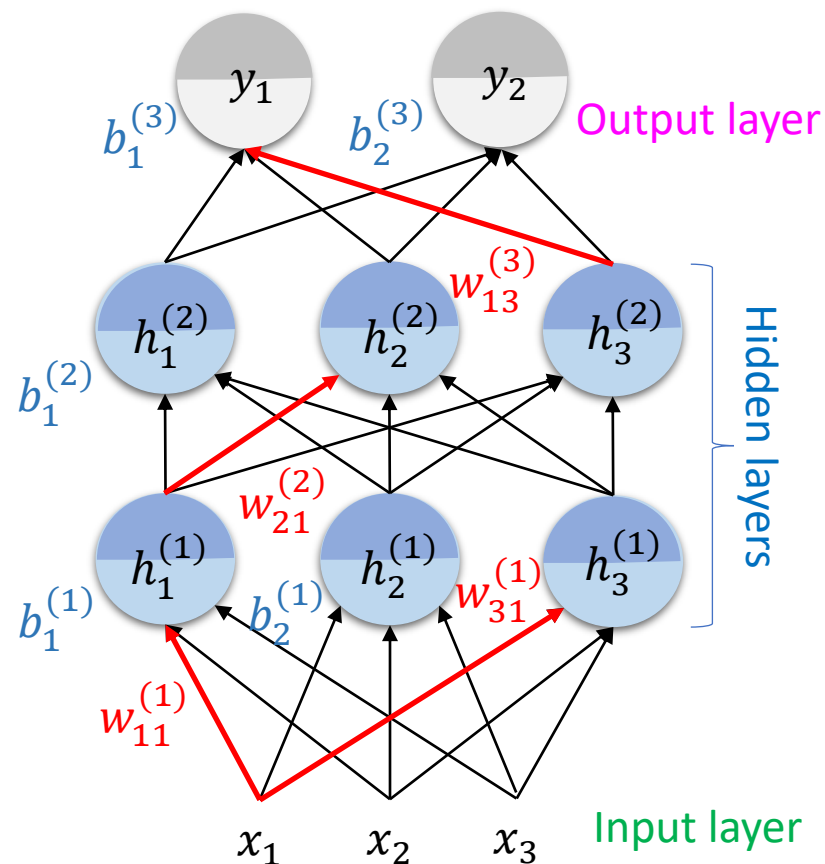
# Computation in the layers of FNN

- Each link (arrow) is associated with a weight
  - $w_{ij}^{(k)} \leftarrow k^{\text{th}} \text{ layer, } i^{\text{th}} \text{ output, } j^{\text{th}} \text{ input}$
- Each circular node has its own bias (not shown)
  - $b_i^{(k)} \leftarrow k^{\text{th}} \text{ layer, } i^{\text{th}} \text{ unit}$

- Network computations:**

- $h_i^{(1)} = \sigma^{(1)} \left( \underbrace{\sum_j w_{ij}^{(1)} x_j}_{\text{Preactivation}} + b_i^{(1)} \right)$   
*activation*
- $h_i^{(2)} = \sigma^{(2)} \left( \sum_j w_{ij}^{(2)} h_j^{(1)} + b_i^{(2)} \right)$
- $y_i = \sigma^{(3)} \left( \sum_j w_{ij}^{(3)} h_j^{(2)} + b_i^{(3)} \right)$

- Note that  $\sigma^{(1)}$  and  $\sigma^{(2)}$  can be different because different layers may have different activation functions



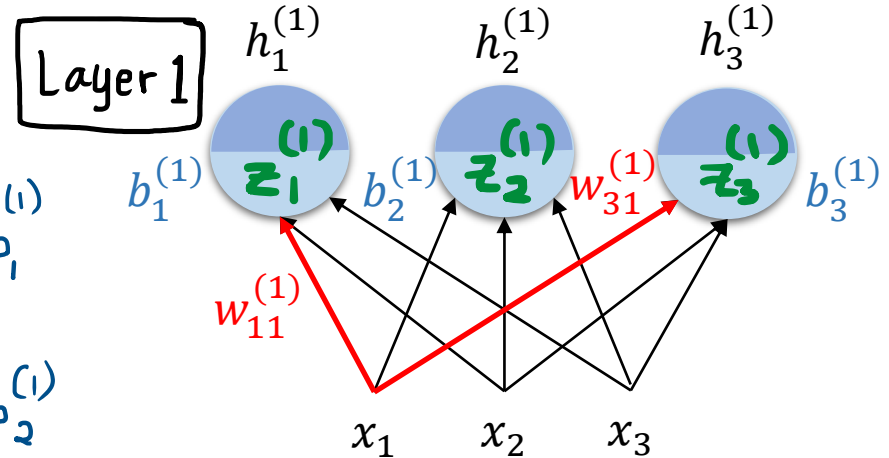
# Computation in the layers of FNN

- Compute preactivation

$$z_1^{(1)} = w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + w_{13}^{(1)} x_3 + b_1^{(1)}$$

$$z_2^{(1)} = w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + w_{23}^{(1)} x_3 + b_2^{(1)}$$

$$z_3^{(1)} = w_{31}^{(1)} x_1 + w_{32}^{(1)} x_2 + w_{33}^{(1)} x_3 + b_3^{(1)}$$



- Compute activation

$$h_1^{(1)} = \sigma^{(1)}(z_1^{(1)})$$

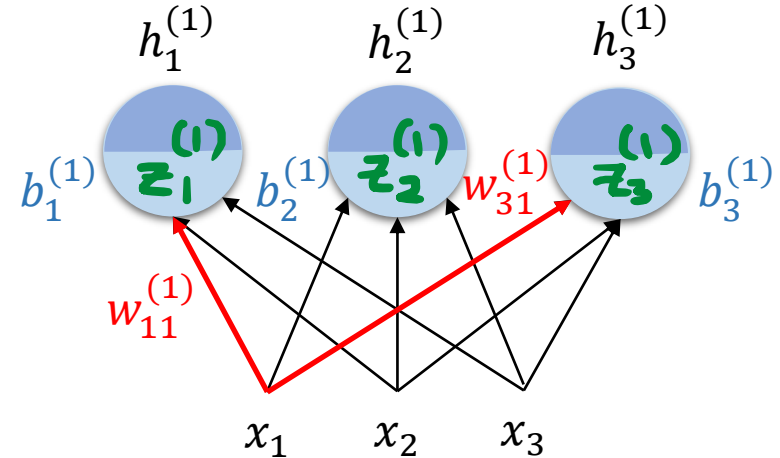
$$h_2^{(1)} = \sigma^{(1)}(z_2^{(1)})$$

$$h_3^{(1)} = \sigma^{(1)}(z_3^{(1)})$$

# Computation in the layers of FNN

- Compute preactivation

$$\begin{pmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \end{pmatrix} = \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} & w_{13}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} & w_{23}^{(1)} \\ w_{31}^{(1)} & w_{32}^{(1)} & w_{33}^{(1)} \end{bmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} + \begin{pmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \end{pmatrix}$$



$$\tilde{z}^{(1)} = \tilde{W}^{(1)} \tilde{x} + \tilde{b}^{(1)} \quad \leftarrow \text{vectorized}$$

- Compute activation

$$h^{(1)} = \sigma^{(1)}(\tilde{z}^{(1)})$$

$$\begin{pmatrix} h_1^{(1)} \\ h_2^{(1)} \\ h_3^{(1)} \end{pmatrix} = \begin{pmatrix} \sigma^{(1)}(z_1^{(1)}) \\ \sigma^{(1)}(z_2^{(1)}) \\ \sigma^{(1)}(z_3^{(1)}) \end{pmatrix}$$

# Computation in the layers of FNN

- $\mathbf{x} = [x_1 \quad x_2 \quad x_3]^T, \mathbf{y} = [y_1 \quad y_2]^T$

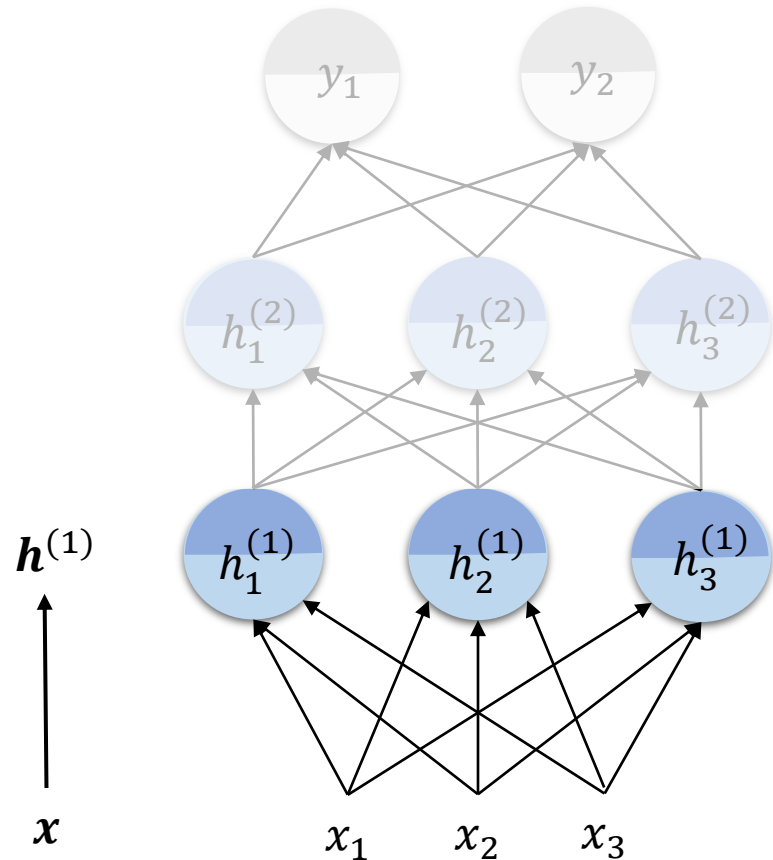
- $\mathbf{h}^{(1)} = [h_1^{(1)} \quad h_2^{(1)} \quad h_3^{(1)}]^T$

- $\mathbf{b}^{(1)} = [b_1^{(1)} \quad b_2^{(1)} \quad b_3^{(1)}]^T$

- $\mathbf{W}^{(1)} = \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} & w_{13}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} & w_{23}^{(1)} \\ w_{31}^{(1)} & w_{32}^{(1)} & w_{33}^{(1)} \end{bmatrix}$

- **Vectorized computation for 1<sup>st</sup> layer**

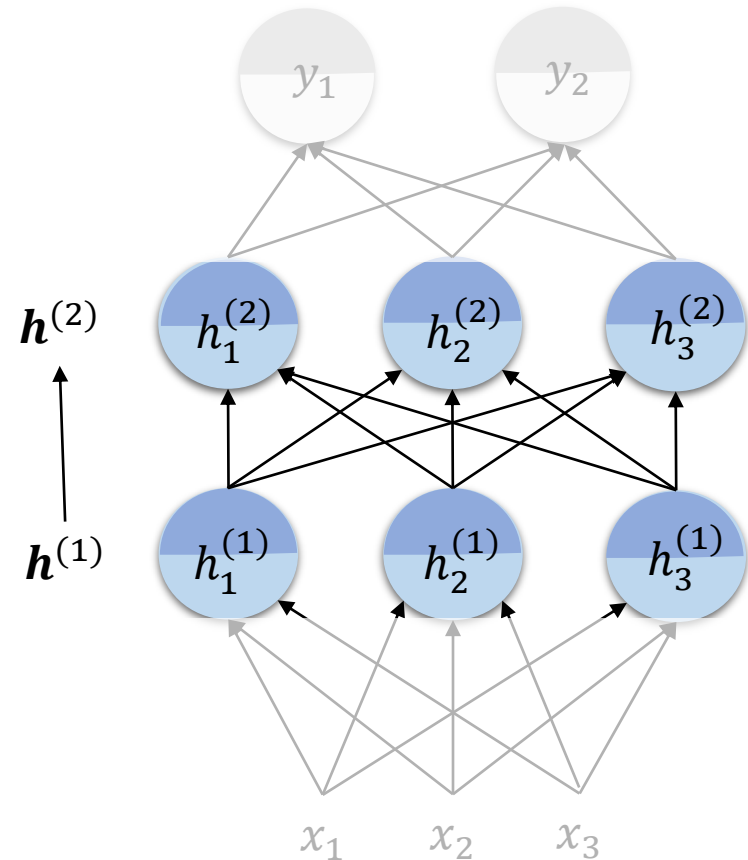
$$\mathbf{h}^{(1)} = \sigma^{(1)}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$$



# Computation in the layers of FNN

- $\mathbf{h}^{(1)} = [h_1^{(1)} \quad h_2^{(1)} \quad h_3^{(1)}]^T$
- $\mathbf{h}^{(2)} = [h_1^{(2)} \quad h_2^{(2)} \quad h_3^{(2)}]^T$
- $\mathbf{b}^{(2)} = [b_1^{(2)} \quad b_2^{(2)} \quad b_3^{(2)}]^T$
- $\mathbf{W}^{(2)} = \begin{bmatrix} w_{11}^{(2)} & w_{12}^{(2)} & w_{13}^{(2)} \\ w_{21}^{(2)} & w_{22}^{(2)} & w_{23}^{(2)} \\ w_{31}^{(2)} & w_{32}^{(2)} & w_{33}^{(2)} \end{bmatrix}$
- **Vectorized computation for 2<sup>nd</sup> layer**

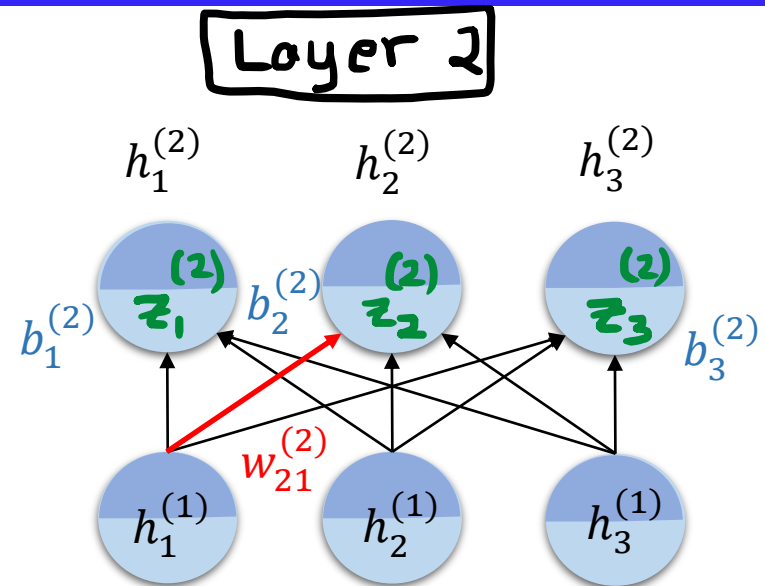
$$\mathbf{h}^{(2)} = \sigma^{(2)}(\mathbf{W}^{(2)}\mathbf{h}^{(1)} + \mathbf{b}^{(2)})$$



# Computation in the layers of FNN

- Compute preactivation

$$\begin{pmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{pmatrix} = \begin{bmatrix} w_{11}^{(2)} & w_{12}^{(2)} & w_{13}^{(2)} \\ w_{21}^{(2)} & w_{22}^{(2)} & w_{23}^{(2)} \\ w_{31}^{(2)} & w_{32}^{(2)} & w_{33}^{(2)} \end{bmatrix} \begin{pmatrix} h_1^{(1)} \\ h_2^{(1)} \\ h_3^{(1)} \end{pmatrix} + \begin{pmatrix} b_1^{(2)} \\ b_2^{(2)} \\ b_3^{(2)} \end{pmatrix}$$



$$\tilde{z}^{(2)} = \tilde{W}^{(2)} \tilde{h}^{(1)} + \tilde{b}^{(2)} \quad \leftarrow \text{vectorized}$$

- Compute activation

$$h^{(2)} = \sigma^{(2)}(\tilde{z}^{(2)})$$

$$\begin{pmatrix} h_1^{(2)} \\ h_2^{(2)} \\ h_3^{(2)} \end{pmatrix} = \begin{pmatrix} \sigma^{(2)}(z_1^{(2)}) \\ \sigma^{(2)}(z_2^{(2)}) \\ \sigma^{(2)}(z_3^{(2)}) \end{pmatrix}$$

# Vectorized computation in FNN

- $\mathbf{x} = [x_1 \quad x_2 \quad x_3]^T, \mathbf{y} = [y_1 \quad y_2]^T$

- $\mathbf{h}^{(2)} = [h_1^{(2)} \quad h_2^{(2)} \quad h_3^{(2)}]^T$

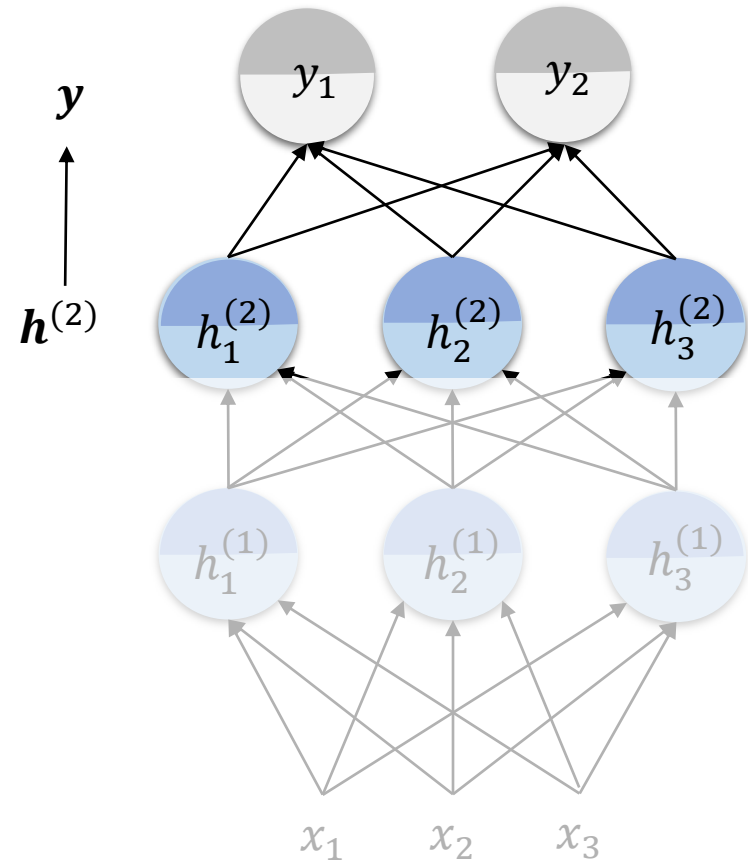
- $\mathbf{b}^{(3)} = [b_1^{(3)} \quad b_2^{(3)} \quad b_3^{(3)}]^T$

- $\mathbf{W}^{(3)} = \begin{bmatrix} w_{11}^{(3)} & w_{12}^{(3)} & w_{13}^{(3)} \\ w_{21}^{(3)} & w_{22}^{(3)} & w_{23}^{(3)} \\ w_{31}^{(3)} & w_{32}^{(3)} & w_{33}^{(3)} \end{bmatrix}$

- **Vectorized computation for 3<sup>rd</sup> layer**

$$\mathbf{y} = \sigma^{(3)}(\mathbf{W}^{(3)}\mathbf{h}^{(2)} + \mathbf{b}^{(3)})$$

- For the output layer, the type of activation function will depend upon the task at hand. It could be linear for regression or sigmoid for classification.

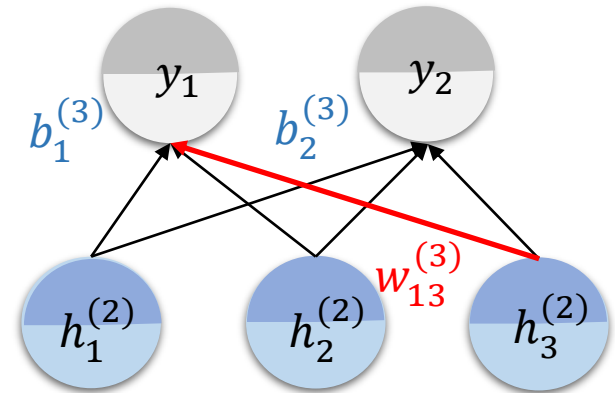




# Computation in the layers of FNN

- Compute preactivation

$$\begin{pmatrix} \tilde{z}_1^{(3)} \\ \tilde{z}_2^{(3)} \end{pmatrix} = \begin{bmatrix} w_{11}^{(3)} & w_{12}^{(3)} & w_{13}^{(3)} \\ w_{21}^{(3)} & w_{22}^{(3)} & w_{23}^{(3)} \end{bmatrix} \begin{pmatrix} h_1^{(2)} \\ h_2^{(2)} \\ h_3^{(2)} \end{pmatrix} + \begin{pmatrix} b_1^{(3)} \\ b_2^{(3)} \end{pmatrix}$$



$$\tilde{\mathbf{z}}^{(3)} = \mathbf{W}^{(3)} \mathbf{h}^{(2)} + \mathbf{b}^{(3)} \quad \leftarrow \text{vectorized}$$

- Compute activation

$$\mathbf{y} = \sigma^{(3)}(\tilde{\mathbf{z}}^{(3)})$$

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} \sigma^{(3)}(\tilde{z}_1^{(3)}) \\ \sigma^{(3)}(\tilde{z}_2^{(3)}) \end{pmatrix}$$

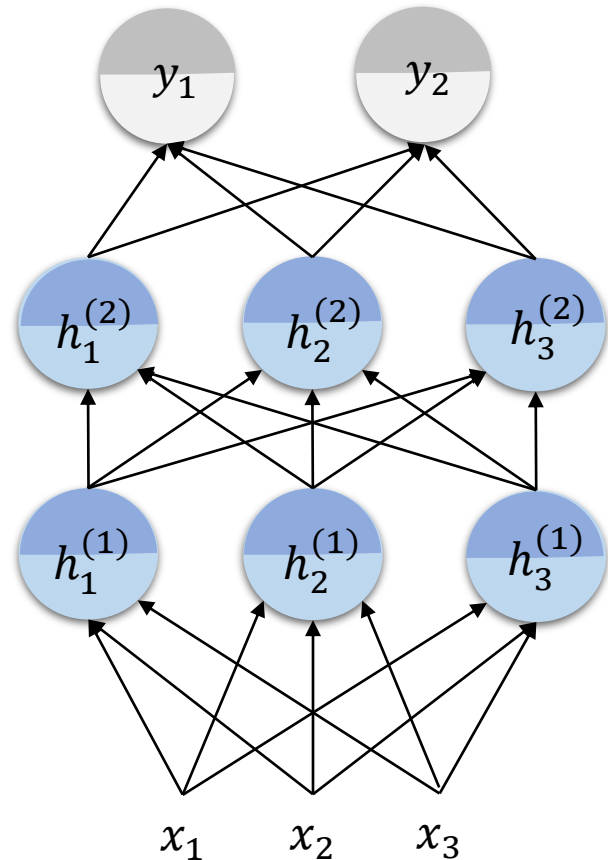
# Vectorized computation in FNN

- Compact vectorized representation

$$\mathbf{h}^{(1)} = \sigma^{(1)}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$$

$$\mathbf{h}^{(2)} = \sigma^{(2)}(\mathbf{W}^{(2)}\mathbf{h}^{(1)} + \mathbf{b}^{(2)})$$

$$\mathbf{y} = \sigma^{(3)}(\mathbf{W}^{(3)}\mathbf{h}^{(2)} + \mathbf{b}^{(3)})$$



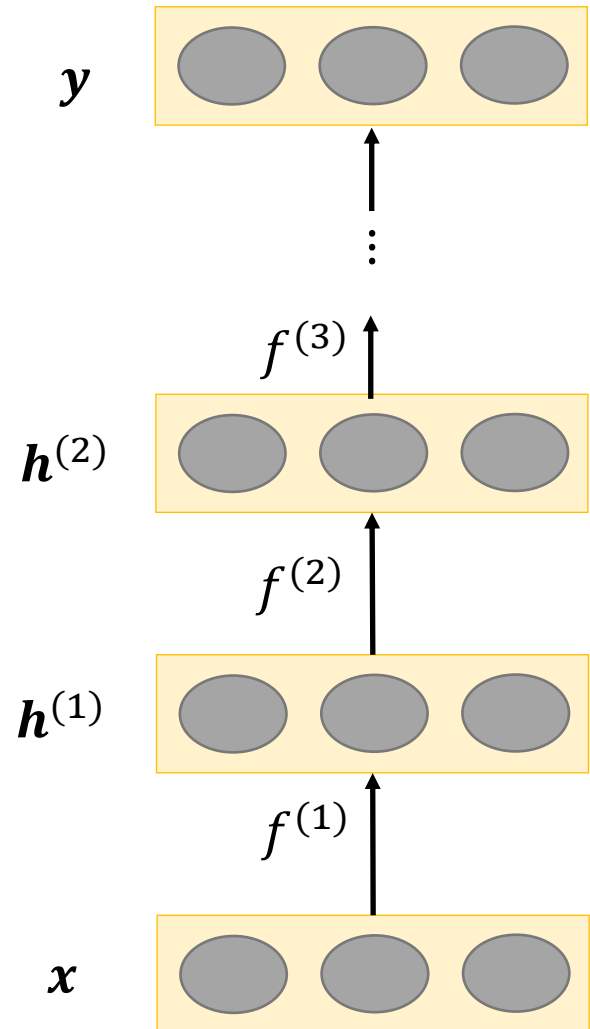
# FNN as a function of functions

- How does the output  $\mathbf{y}$  related to the input  $\mathbf{x}$
- Each layer computes a function, so the network computes a **composition of functions**:

$$\begin{aligned}\mathbf{h}^{(1)} &= f^{(1)}(\mathbf{x}) \\ \mathbf{h}^{(2)} &= f^{(2)}(\mathbf{h}^{(1)}) \\ &\vdots \\ \mathbf{y} &= f^{(L)}(\mathbf{h}^{(L-1)})\end{aligned}$$

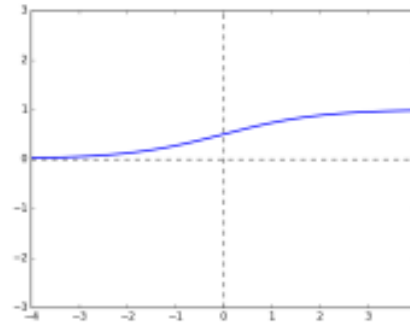
- Or more simply

$$\mathbf{y} = f^{(L)} \left( f^{(L-1)} \left( \dots f^{(1)}(\mathbf{x}) \right) \right)$$



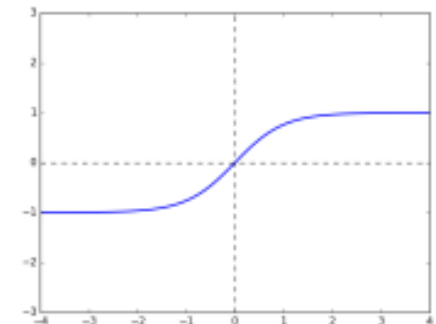
# Activation functions

- We are free to use any kind of differentiable activation function we like at each layer
- A common choice was to use a sigmoid (logistic) function, which can be seen as a smooth approximation to the step function used in a perceptron



**Logistic**

$$y = \frac{1}{1 + e^{-z}}$$



**Hyperbolic Tangent (tanh)**

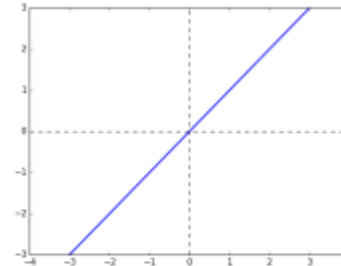
$$y = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

- However, the **logistic** function saturates at 1 for large positive inputs, and at 0 for large negative inputs
- The **tanh function** has a similar shape, but saturates at -1 and +1. In these regimes, the gradient of the output w.r.t. the input will be close to zero, so any gradient signal from higher layers will “vanish”, as we will discuss later
- Tanh is better than sigmoid because it is zero-centred

# Activation functions

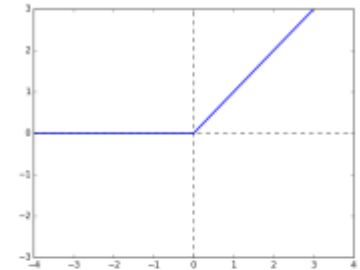
- Non-saturating activation functions
- The most common one is rectified linear unit or **ReLU**

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



**Linear**

$$y = z$$



**Rectified Linear Unit  
(ReLU)**

$$y = \max(0, z)$$

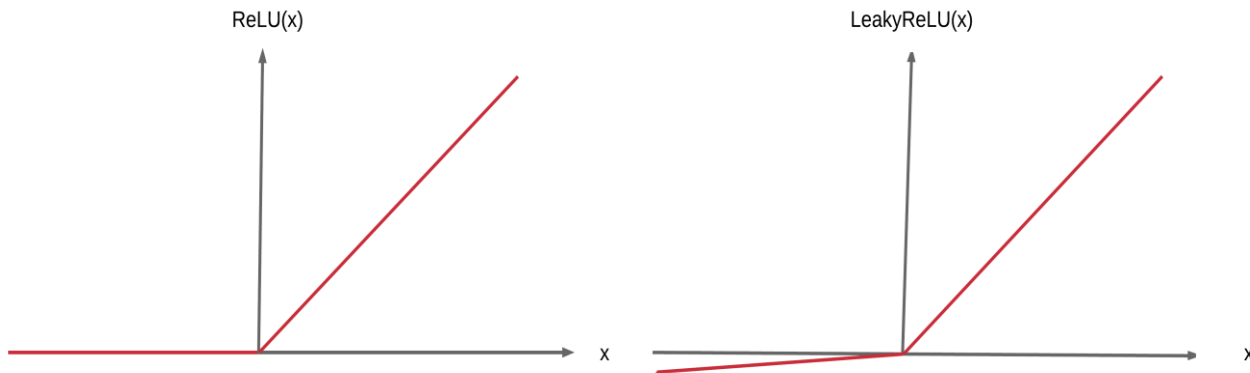
- The ReLU function simply “turns off” negative inputs, and passes positive inputs unchanged. This turns out to have a gradient of 1, at least for positive inputs; this helps avoid vanishing gradients
- Unfortunately, for negative inputs, the gradient of ReLU units is 0, so the unit will never get any feedback signal to help it escape from this parameter setting; this is called the “**dying ReLU**” problem
- ReLU is less computationally expensive than tanh and sigmoid because it involves simpler mathematical operations

# Activation functions

- One simple fix is to use the **leaky ReLU**. This is defined as

$$\text{LReLU}(z; \alpha) = \max(\alpha z, z), \quad 0 < \alpha < 1$$

- The slope of this function is 1 for positive inputs, and  $\alpha$  for negative inputs, thus ensuring there is some signal passed back to earlier layers, even when the input is negative.
- Leaky ReLU is non-smooth, it has a kink at  $z = 0$

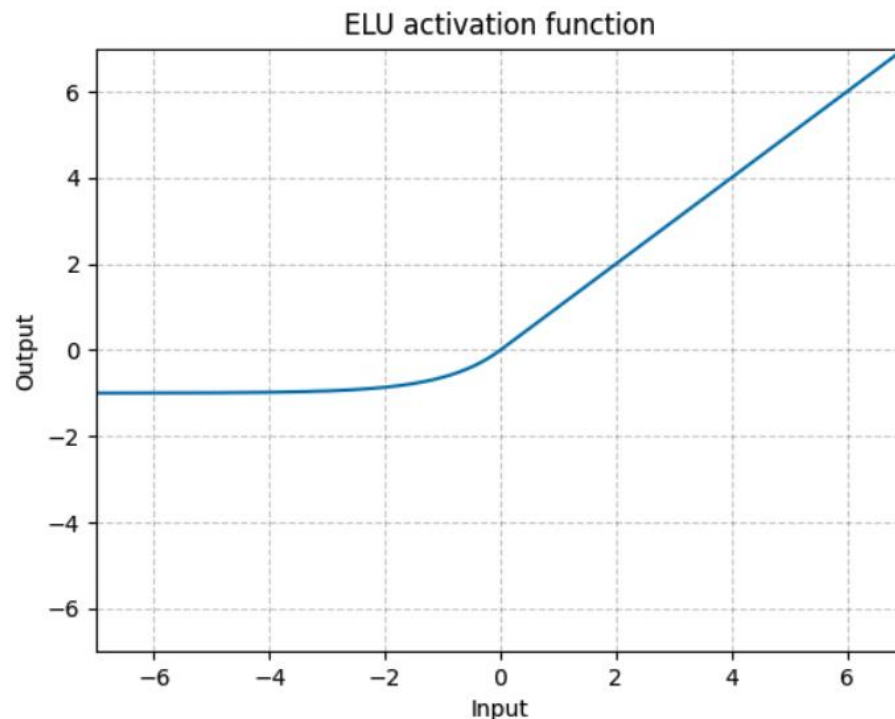


# Activation functions

- Another popular choice is the **ELU**

$$ELU(z; \alpha) = \begin{cases} \alpha(e^{\{z\}} - 1) & \text{if } z \leq 0 \\ z & \text{if } z > 0 \end{cases}$$

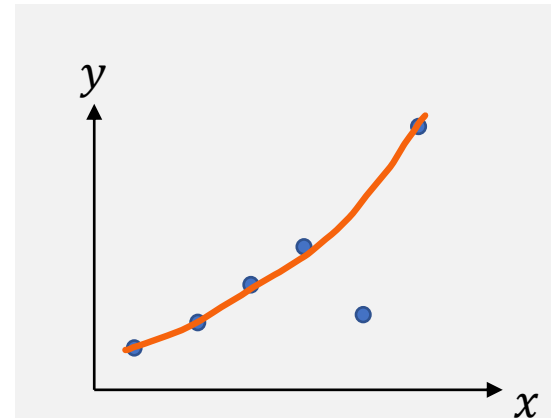
- This has the advantage over leaky ReLU of being a smooth function



# Feature learning

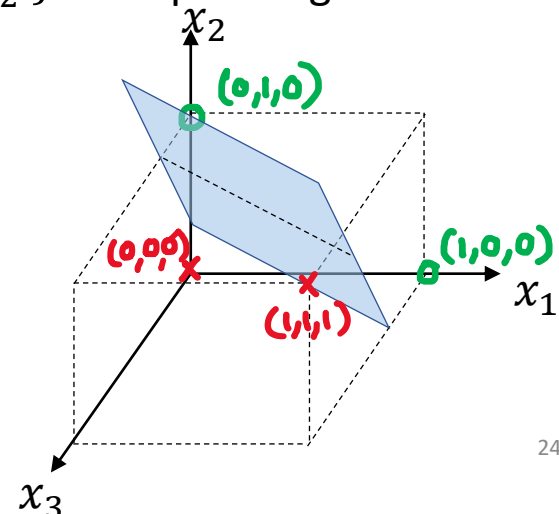
- Recall nonlinear feature maps in linear regression and classification
- For instance, the feature mapping  $\phi(x) = \{1, x, x^2, x^3\}$  can represent third-degree polynomials in regression.

$$y = w_1x + w_2x^2 + w_3x^3 + b$$



- In another instance, we used  $\phi(x_1, x_2) = \{x_1, x_2, x_1x_2\}$  for separating XOR function

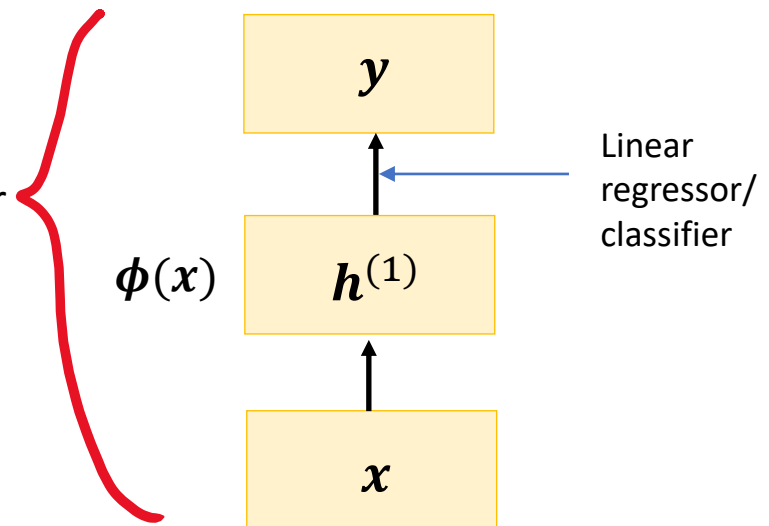
Input $x_1$	Input $x_2$	Input $x_3 = x_1x_2$	Output $y$	
0	0	0	0	✗
1	0	0	1	○
0	1	0	1	○
1	1	1	0	✗





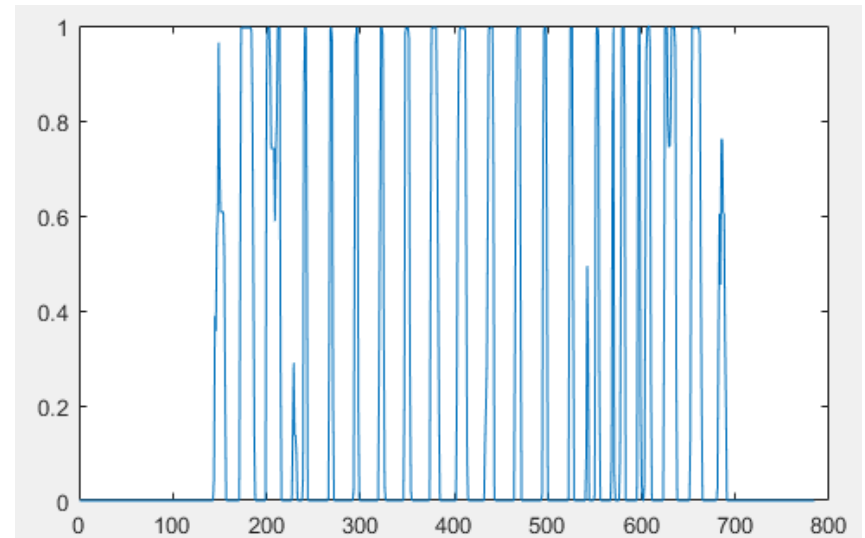
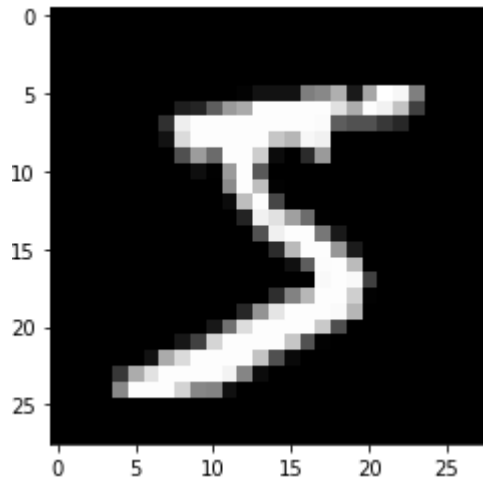
# Feature learning

- Recall nonlinear feature maps in linear regression and classification
- For instance, the feature mapping  $\phi(x) = \{1, x, x^2, x^3\}$  can represent third-degree polynomials in regression.
- In another instance, we used  $\phi(x_1, x_2) = \{x_1, x_2, x_1x_2\}$  for separating XOR function
- Neural nets** can be viewed as a way of learning nonlinear features
- For example, the hidden layer can be thought of as a feature map  $\phi(x)$ , and the output layer weights can be thought of as a linear model using those features



# Feature learning example using MNIST

A multi-class classification task with 10 categories, one for each digit class

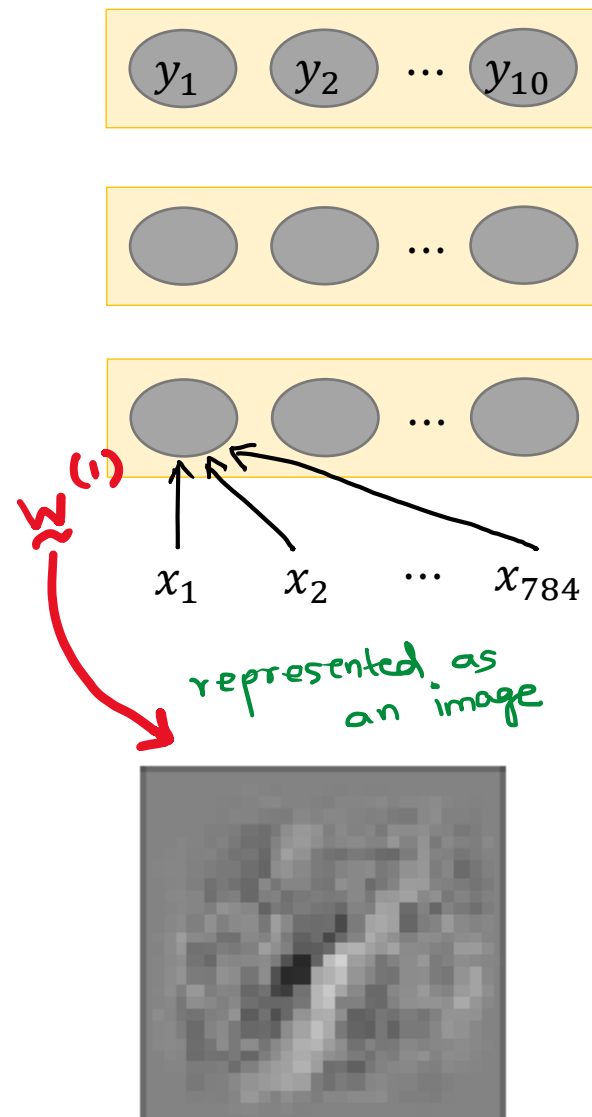


$28 \times 28$   
grayscale  
image

$\mathcal{X}$   
 $784 \times 1$   
input  
vector

# Feature learning example using MNIST

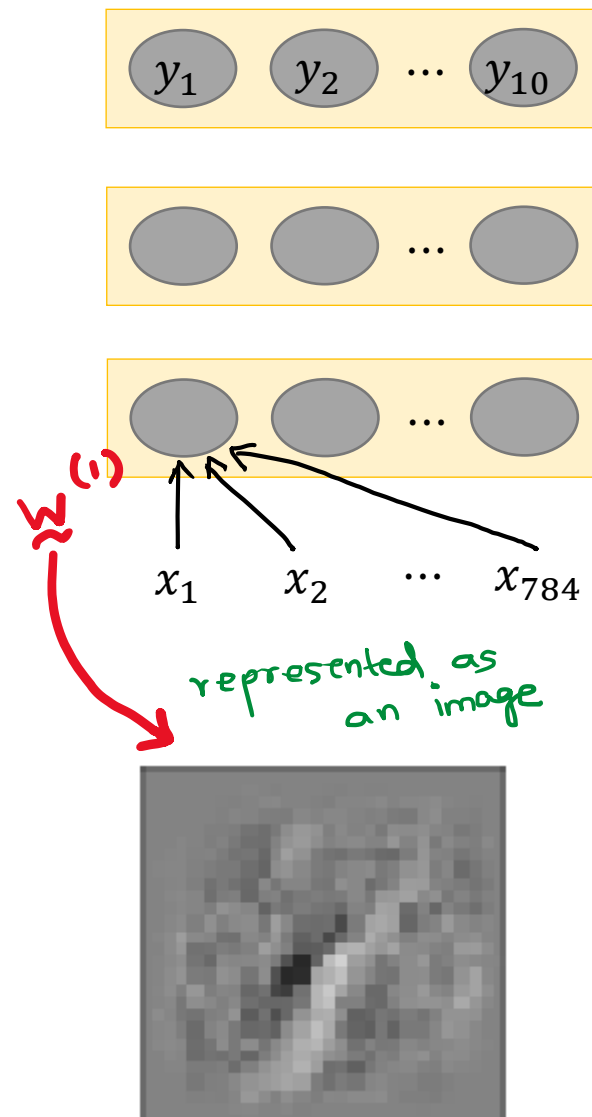
- Suppose we train an MLP with two hidden layers
- Each hidden unit receives inputs from each of the pixels and computes  $\sigma(\mathbf{w}_i^T \mathbf{x})$
- Let's see what the first layer of hidden units computes by visualizing the weights
- The weights  $\mathbf{w}^{(i)} \in \mathbb{R}^{784 \times 1}$  feeding into each hidden unit of the first layer can be reshaped into a  $28 \times 28$  matrix and displayed as images
- Weight image: gray = 0, white  $\rightarrow +$ , black  $\rightarrow -$



# Feature learning example using MNIST

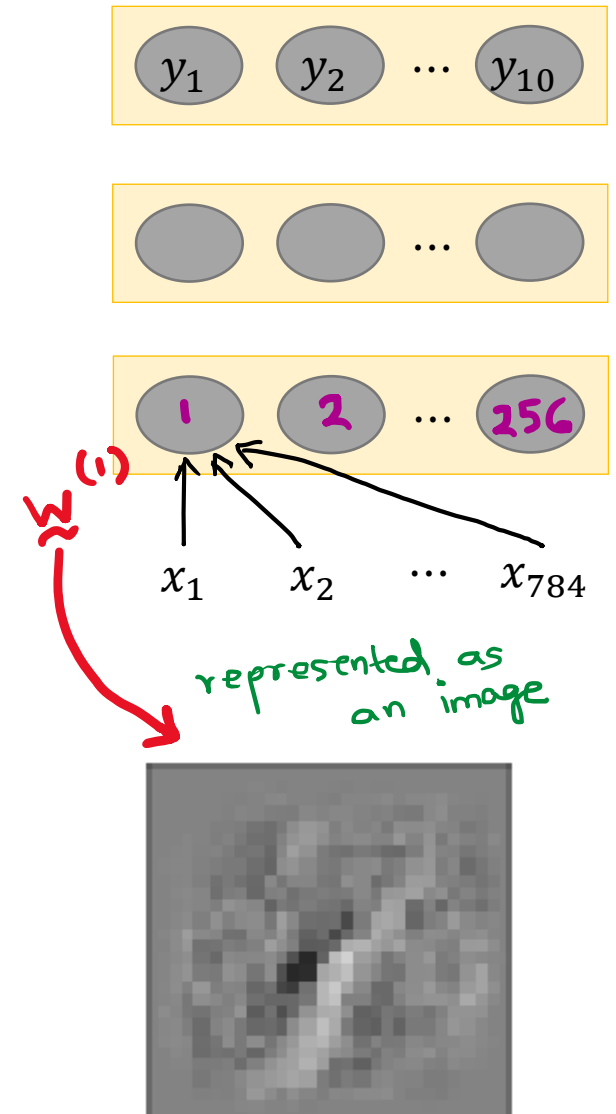
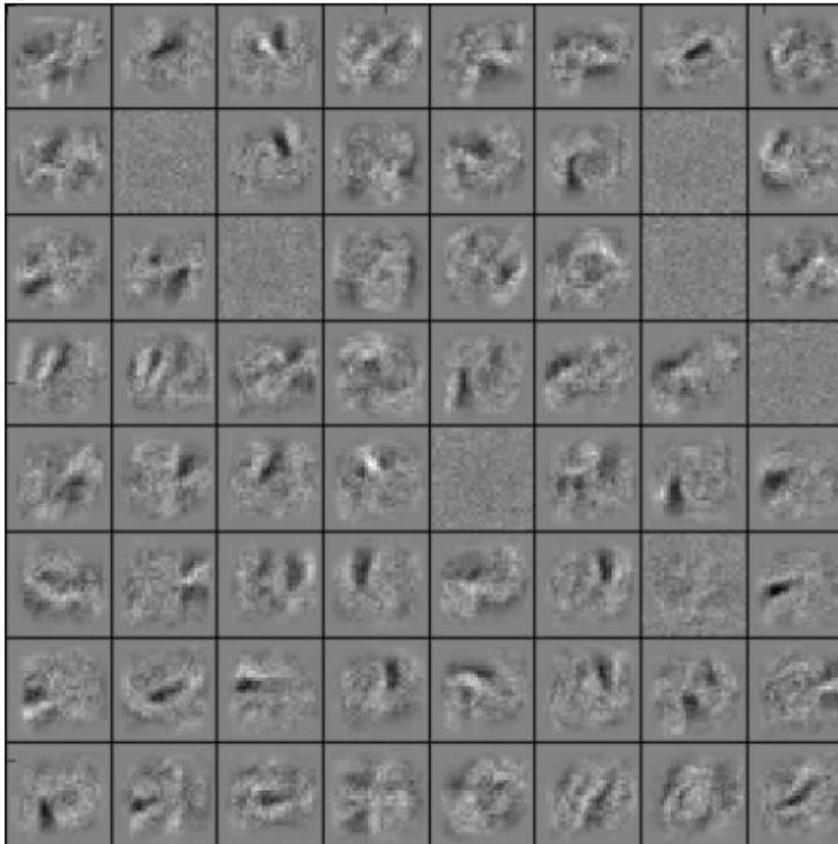
## What do the weights tell us?

- In the weight visualization, positive values are lighter, and negative values are darker
- Each hidden unit computes the dot product of these weight vectors with the input image, and then passes the result through the activation function
- The weight images are mostly found to correspond to some **oriented edges**
- Thus the weights carry useful information about locations and orientations of strokes



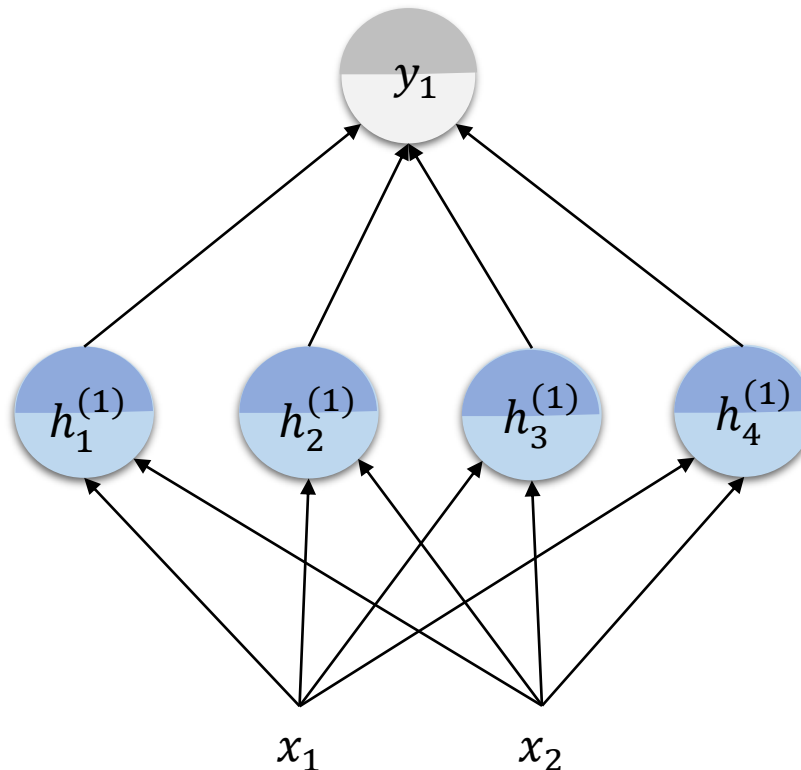
# Feature learning example using MNIST

- There are 256 units in the first hidden layer and hence there are 256 images of weight vectors. Here are some more of images of the weight vectors for the 1<sup>st</sup> layer



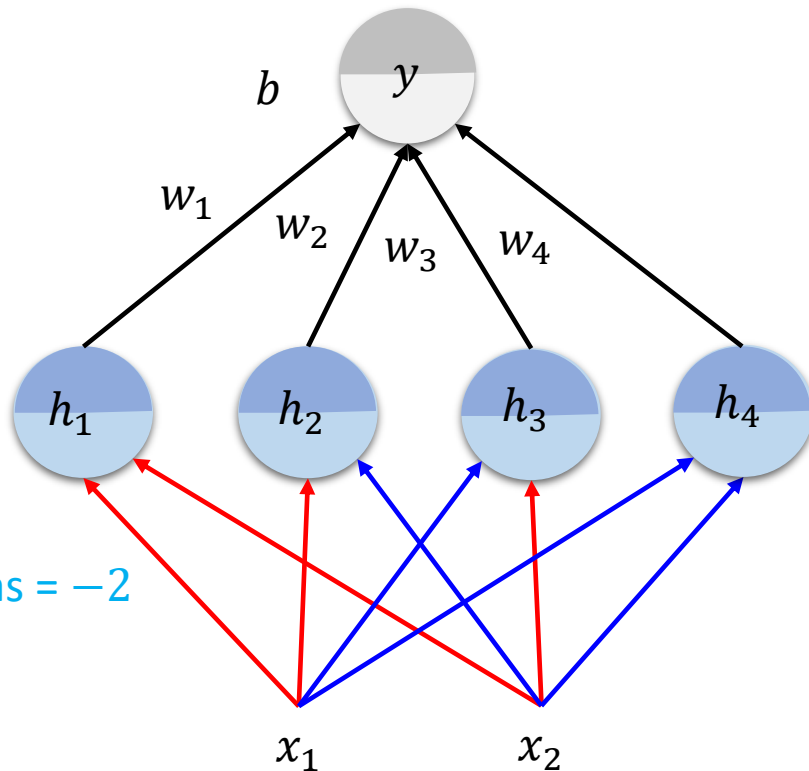
# FNNs are very expressive!

- Multilayer feed-forward neural nets with *nonlinear* activation functions are **universal approximators**: they can approximate any function arbitrarily well
- For example, FNN with **just one hidden layer can represent the any Boolean function** (or logic gates)
- Let's try this network for representing the XOR gate



# FNNs are very expressive!

**Note:** Here we consider hard threshold activation functions for this example



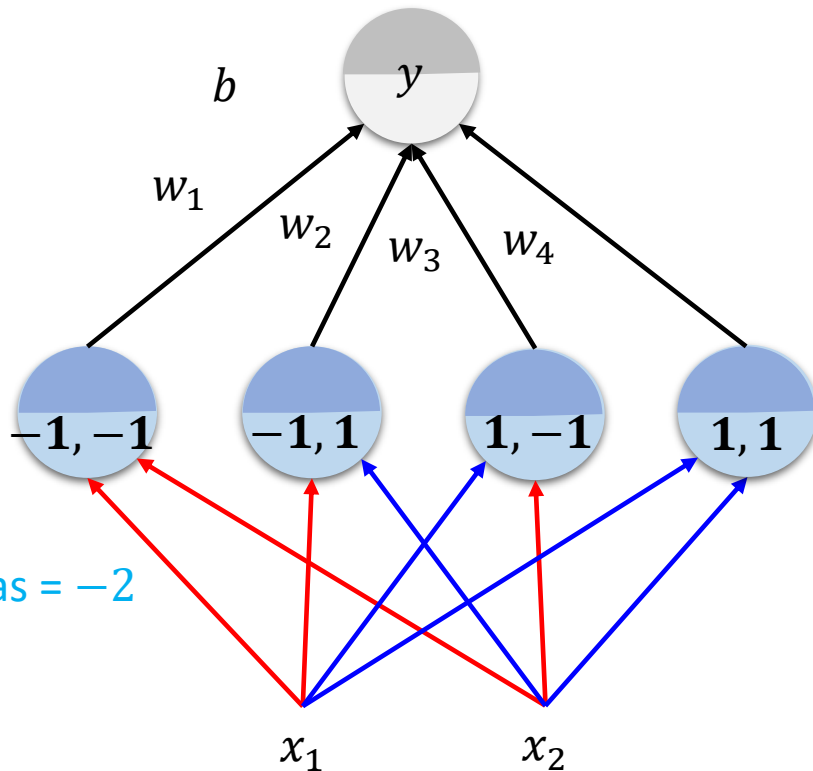
Red edge indicates  $w = -1$

Blue edge indicates  $w = +1$

- For this discussion, we will assume True = +1 and False = -1
- We consider two inputs and four perceptrons in hidden layer
- The weights of the first layer are manually fixed for illustration
- The bias for each perceptron is set to -2
- The weights  $w_i$ 's and bias for the output perceptron are to be learned (in this case)
- **Question:** Can we find  $\{w_1, w_2, w_3, w_4, b\}$  such that the FNN can represent any Boolean function?

# FNNs are very expressive!

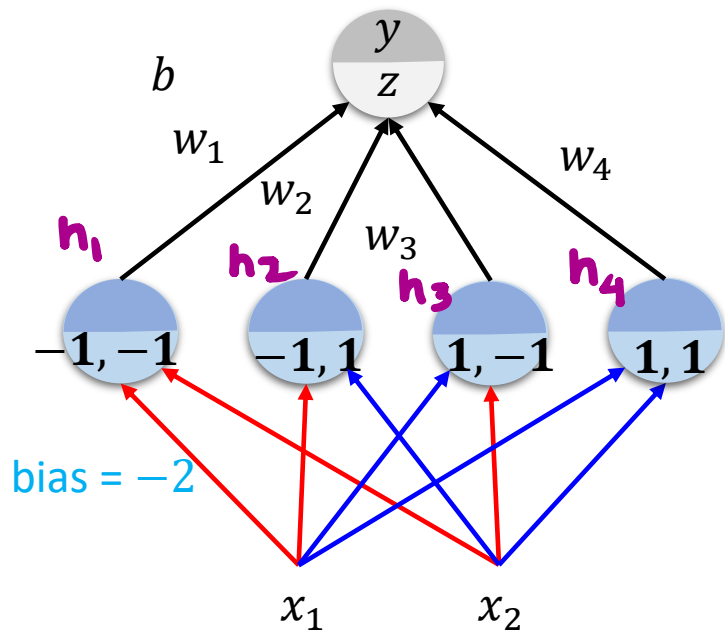
**Note:** Here we consider hard threshold activation functions for this example



- For this discussion, we will assume True =  $+1$  and False =  $-1$
- **Question:** Can we find  $\{w_1, w_2, w_3, w_4, b\}$  such that the FNN can represent any Boolean function?
- Each hidden unit fires for a specific input (no two perceptron fire for the same input)
- The first hidden perceptron fires for inputs {False, False}  $\rightarrow \{-1, -1\}$
- The second hidden perceptron fires for inputs {False, True}  $\rightarrow \{-1, 1\}$



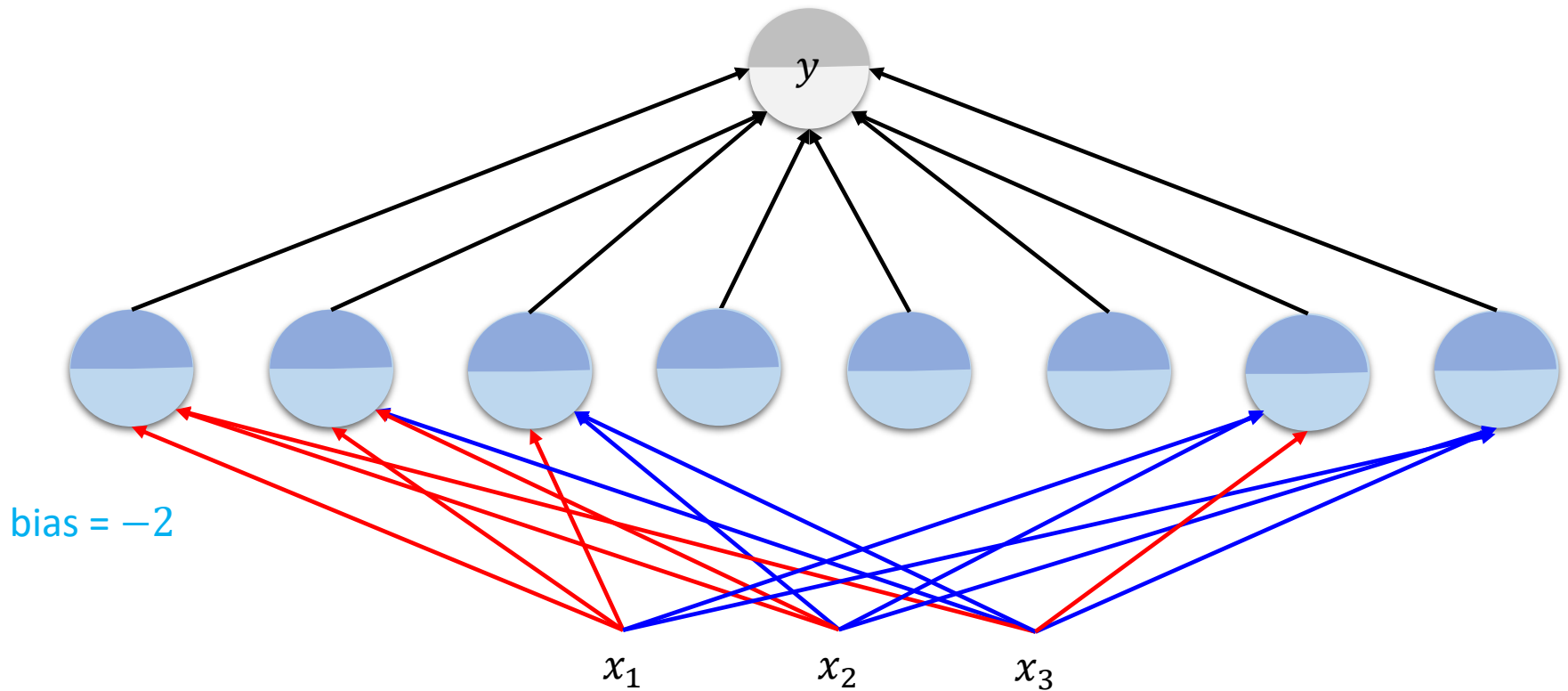
# FNNs are very expressive!



- Can we adjust the weights to get whatever output we need?
- **Say we want to represent XOR**
- $b = -1$
- $w_1 = 0, w_2 = 1, w_3 = 1, w_4 = 0$

$x_1$	$x_2$	XOR	$h_1$	$h_2$	$h_3$	$h_4$	$z = \sum w_i h_i + b$	$y = \sigma(z)$
-1	-1	0	1	0	0	0	$w_1 - 1$	
-1	1	1	0	1	0	0	$w_2 - 1$	
1	-1	1	0	0	1	0	$w_3 - 1$	
1	1	0	0	0	0	1	$w_4 - 1$	

# FNNs are very expressive!



Red edge indicates  $w = +1$

Blue edge indicates  $w = -1$

For  $D$  inputs,  $2^D$  units needed in hidden layer

# Shallow vs Deep neural networks

## ➤ Shallow neural networks

- Shallow neural networks have **one hidden layer** of neurons
- They are can very expressive, **but that may require a large width**
- Having an extremely wide neural network can be computationally very costly

## ➤ Deep neural networks

- Deep neural networks have **more than one hidden layer** of neurons
- They are **expressive** as well as **compact**, i.e. don't need to be extremely wide
- Deep neural networks are therefore more practically computationally