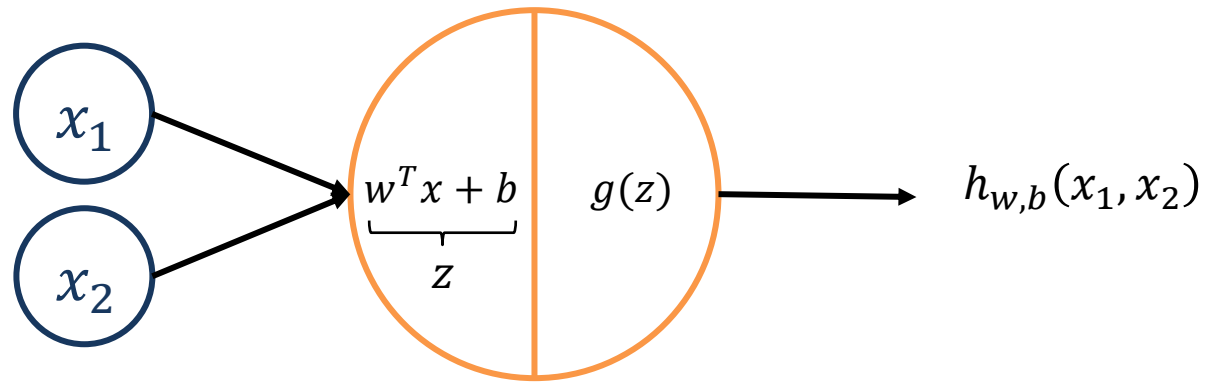


# Neural Networks

Prof. Giuseppe Serra

# Neural Network Representation

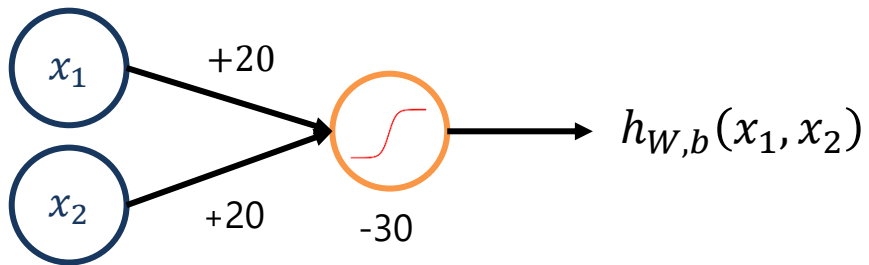
This network means that to compute the value of the decision function, we need to multiply first input  $x_1$  with  $\theta_1$ , second input  $x_2$  with  $\theta_2$ , then add two values and  $b$ , then apply the sigmoid function



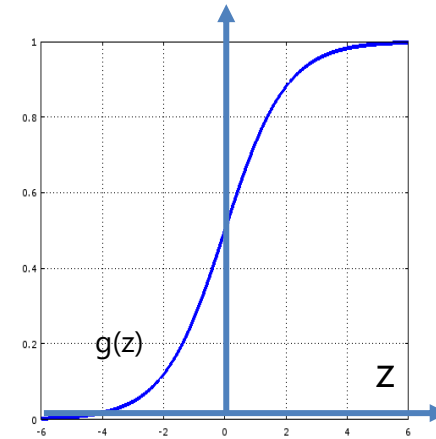
# Simple Neural Networks

Example: AND Logic Port

$$y = x_1 \text{ AND } x_2$$



$$h_{w,b}(x_1, x_2) = g(-30 + 20x_1 + 20x_2)$$

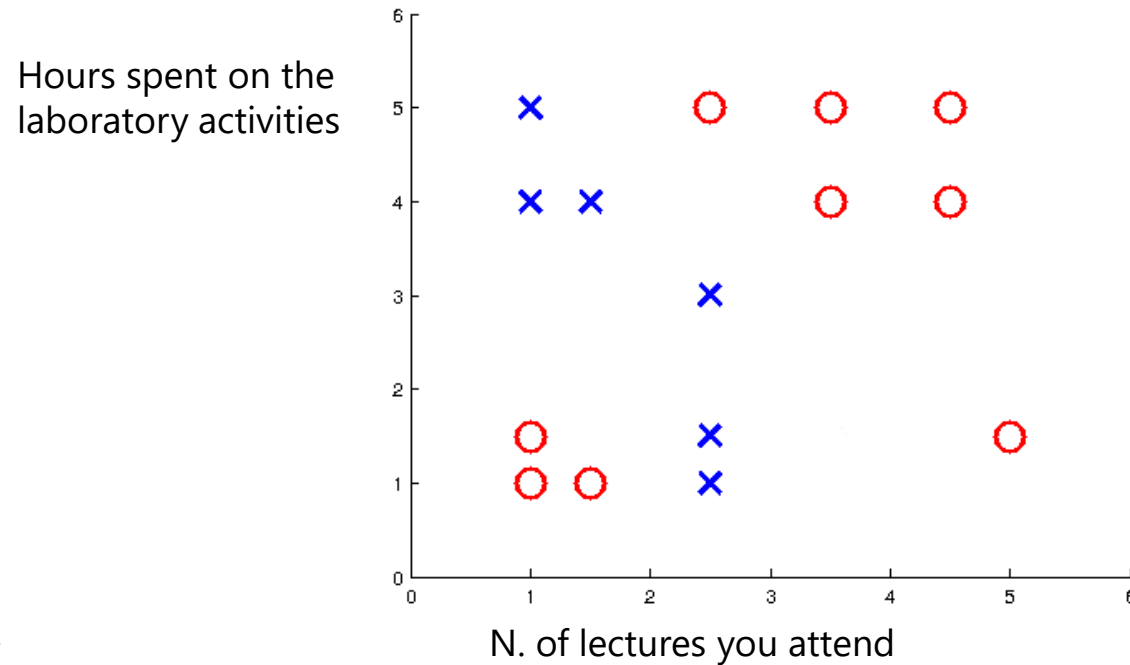


| $x_1$ | $x_2$ | $h_{w,b}(x_1, x_2)$ |
|-------|-------|---------------------|
| 0     | 0     | $g(-30) \approx 0$  |
| 0     | 1     | $g(-10) \approx 0$  |
| 1     | 0     | $g(-10) \approx 0$  |
| 1     | 1     | $g(10) \approx 1$   |

# A Complex Problem

Now consider Susan that has different movies tastes

Logistic Regression Algorithm is not able to solve it.

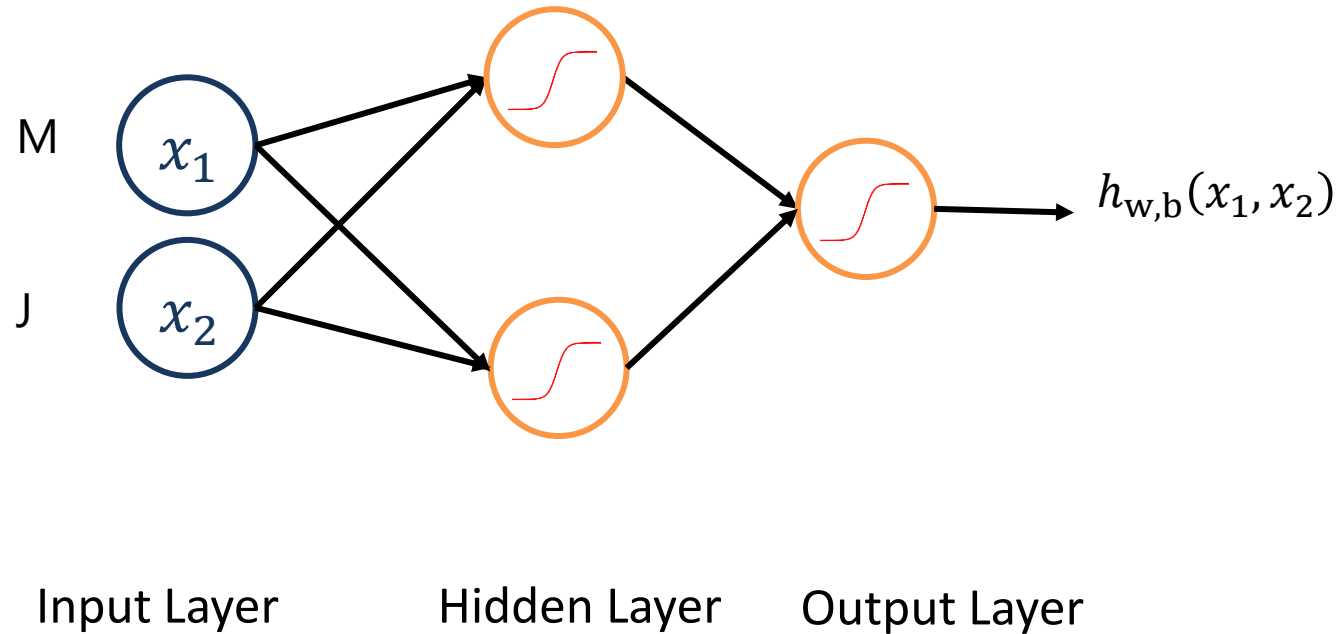
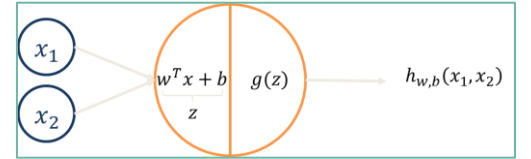


○ corresponds to "Pass"

× corresponds to "Fail"

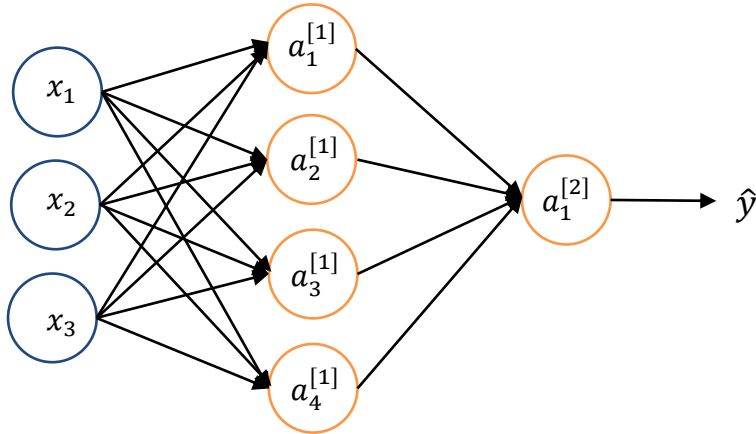
# Neural Network with One Hidden Layer

Let's try to a Neural Network with one hidden layer!



Let's try by yourself: <https://lecture-demo.ira.uka.de/neural-network-demo/?preset=Binary%20Classifier%20for%20XOR>

# Neural Network Representation



$W^{[j]}$  = matrix of weights controlling function mapping from layer  $j - 1$  to layer  $j$   
 $b^{[j]}$  = bias  
 $a^{[j]}$  = activation of vector for layer  $j$   
 $g^{[j]}$  = activation function for layer  $j$

## Given in input X:

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

(4,1) (4,3) (3,1) (4,1) - Dimensionality

$$a^{[1]} = g^{[1]}(z^{[1]})$$

(4,1) (4,1) - Dimensionality

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

(1,1) (1,4) (4,1) (1,1) - Dimensionality

$$a^{[2]} = g^{[2]}(z^{[2]})$$

(1,1) (1,1) - Dimensionality

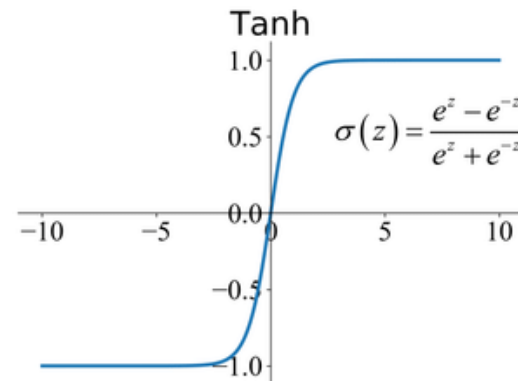
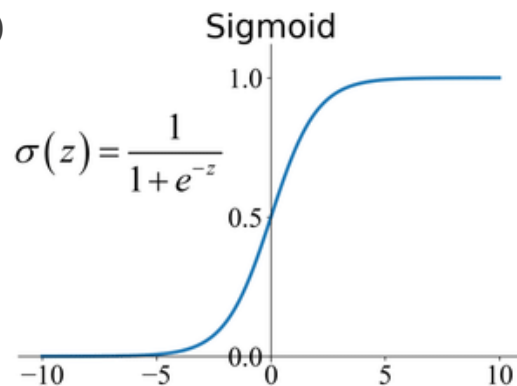
# Activation Functions

## Sigmoid Function:

- it has seen frequent use historically, yet now it is rarely used in the intermediate layers, because it kills the gradient (remember  $\frac{d \sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$ ).
- It is commonly use in the last layer is you want a prediction value between [0,1]

## Tanh Function:

- “...the hyperbolic tangent activation function typically performs better than the logistic sigmoid.” – Page 195, Deep Learning, 2016
- $\frac{d \tanh(z)}{dz} = \dots = 1 - \tanh^2(z)$



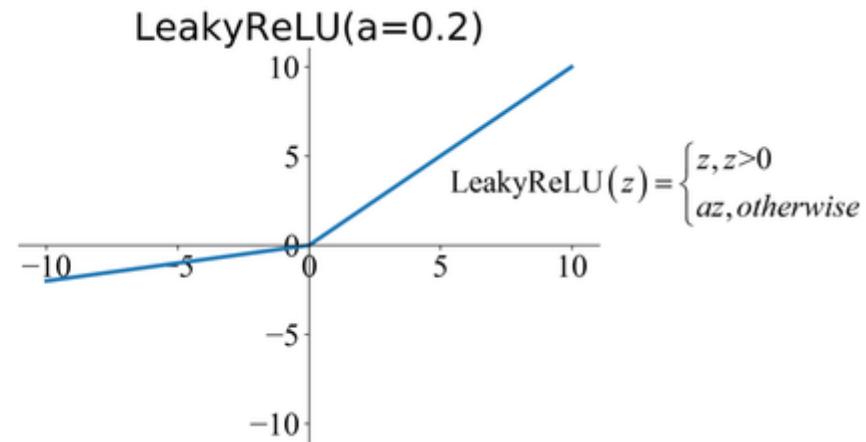
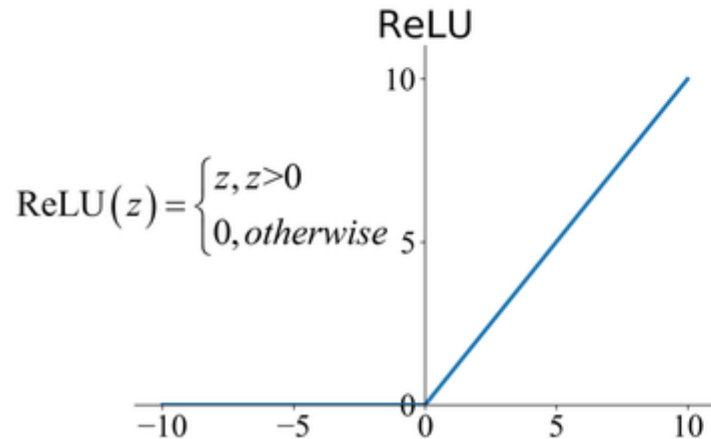
# Activation Functions

Rectified linear unit (ReLU):

- It has become the **default activation function** for many types of neural networks, because a model that uses it is **easier to train** and often achieves better performance. The derivative of the rectified linear function is also **easy to calculate**. **The slope for negative values is 0.0 and the slope for positive values is 1.0. Derivative in zero is not defined** (empirically is not a problem)

Leaky version of a Rectified Linear Unit (LeakyReLU)

- ReLU is affected to the **Dying ReLU Problem**: some **ReLU Neurons essentially die** for all inputs and remain inactive no matter what input is supplied, here no gradient flows (**the gradient of the function becomes zero, the network cannot perform backpropagation and cannot learn**). Empirically performance is similar to the ReLU.





# Why do you need non-linear functions?

## Given in input X:

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

(4,1) (4,3) (3,1) (4,1) - Dimensionality

$$a^{[1]} = g^{[1]}(z^{[1]})$$

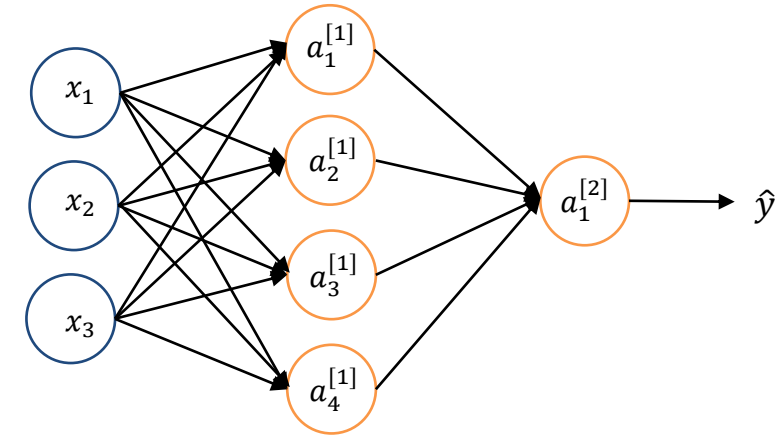
(4,1) (4,1) - Dimensionality

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

(1,1) (1,4) (4,1) (1,1) - Dimensionality

$$a^{[2]} = g^{[2]}(z^{[2]})$$

(1,1) (1,1) - Dimensionality

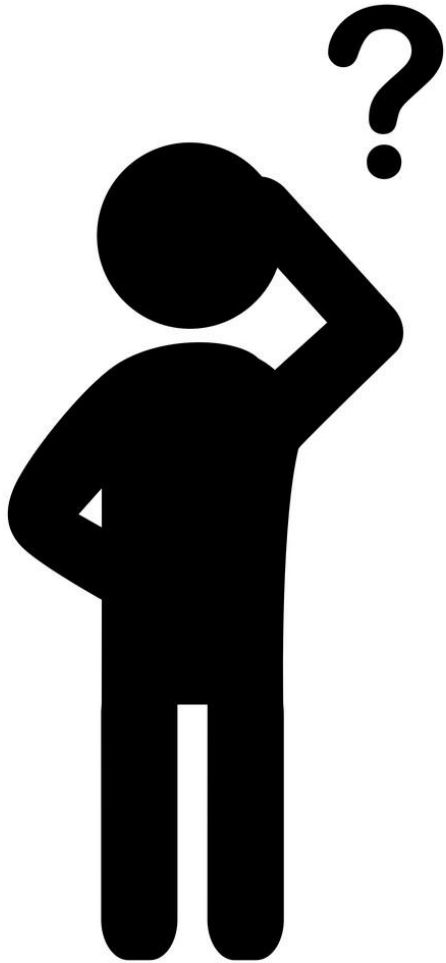


- Let's suppose  $g^{[1]}(z) = g^{[2]}(z) = z$  (identity function or linear activation function) then:
- $a^{[2]} = z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} = W^{[2]}(W^{[1]}x + b^{[1]}) + b^{[2]} = (W^{[2]}W^{[1]})x + (W^{[2]}b^{[1]} + b^{[2]})$
- Let's define  $W' = W^{[2]}W^{[1]}$  and  $b' = (W^{[2]}b^{[1]} + b^{[2]})$
- $a^{[2]} = W'x + b'$  (Thus, if you are using identity activation functions, the output neural networks is linear combination of the input)
- Linear activation is often useless, unless in the last layer of a neural network when you are working in a regression task.
- Let's try again by yourself (use linear activation functions): <https://lecture-demo.ira.uka.de/neural-network-demo/?preset=Binary%20Classifier%20for%20XOR>

# Loss Function

## Example Problem: Will I Pass this Class?

---

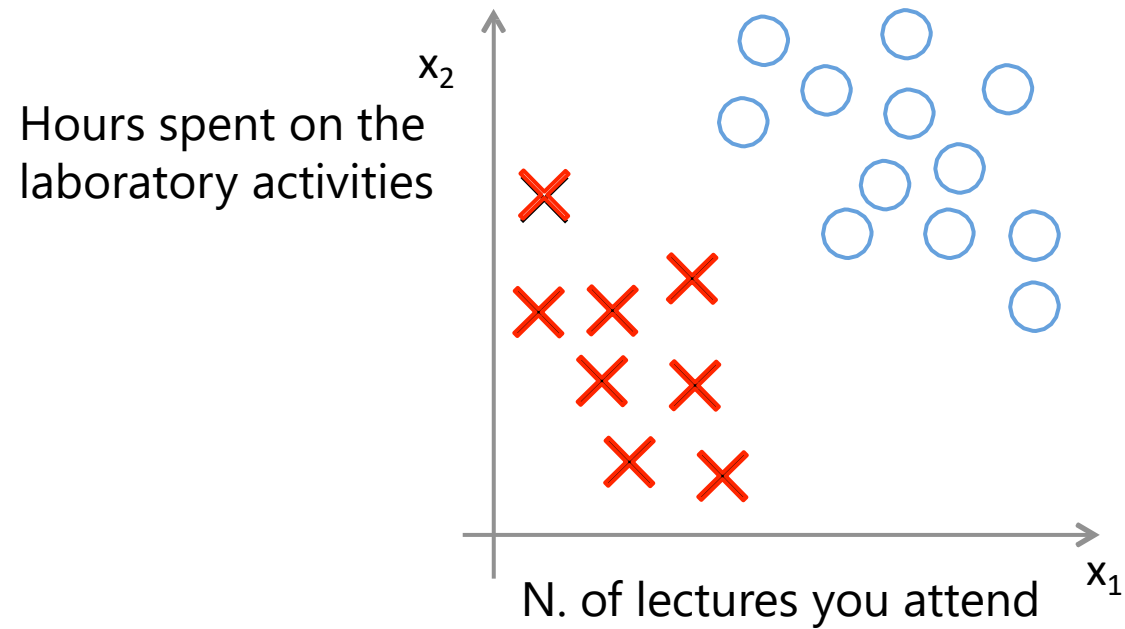


Let's start with a simple two feature model:

$x_1$  = Number of lectures you attend

$x_2$  = Hours spent on the laboratories activities

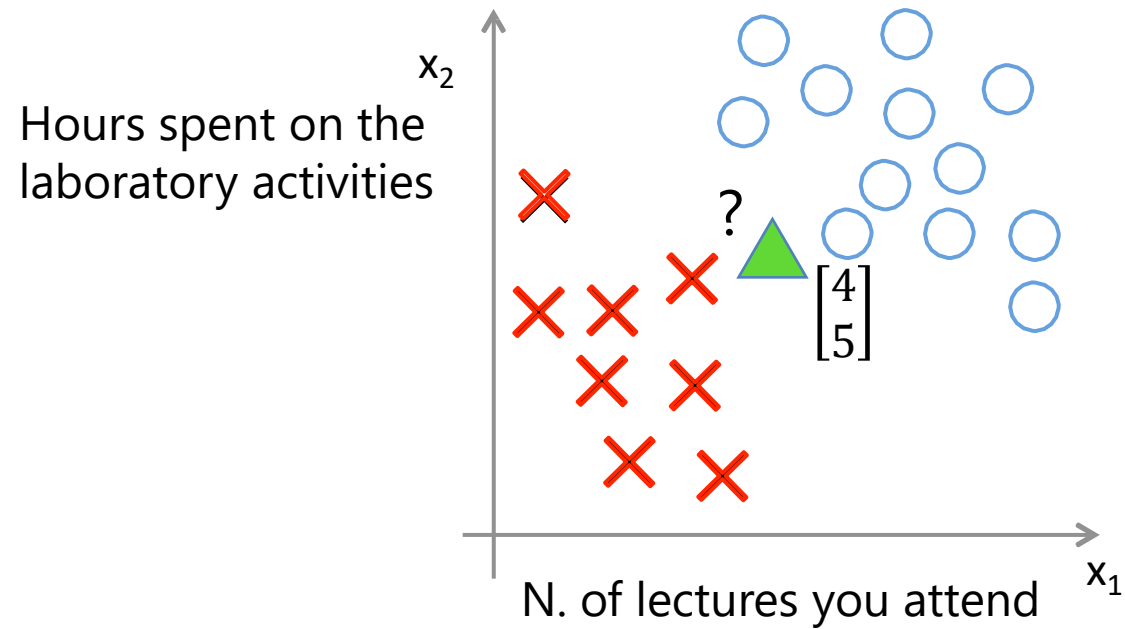
# Example Problem: Will I Pass this Class?



○ corresponds to "Pass"

✗ corresponds to "Fail"

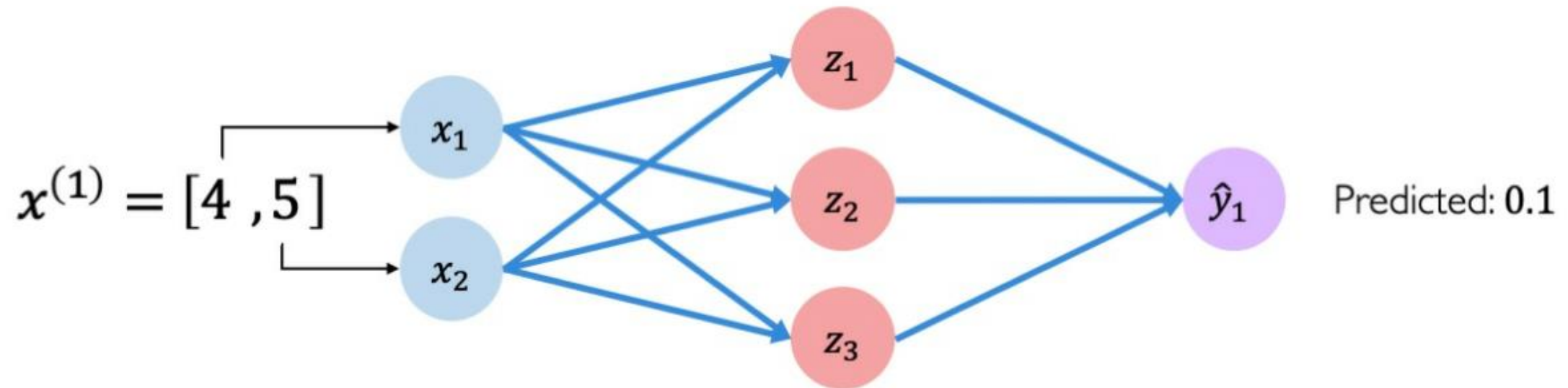
# Example Problem: Will I Pass this Class?



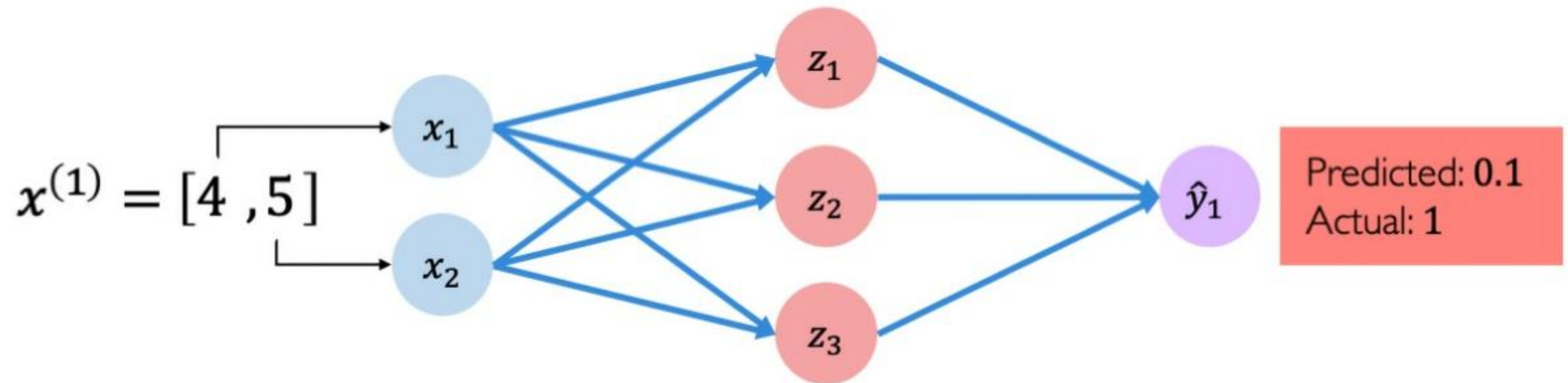
○ corresponds to "Pass"

✗ corresponds to "Fail"

## Example Problem: Will I Pass this Class?

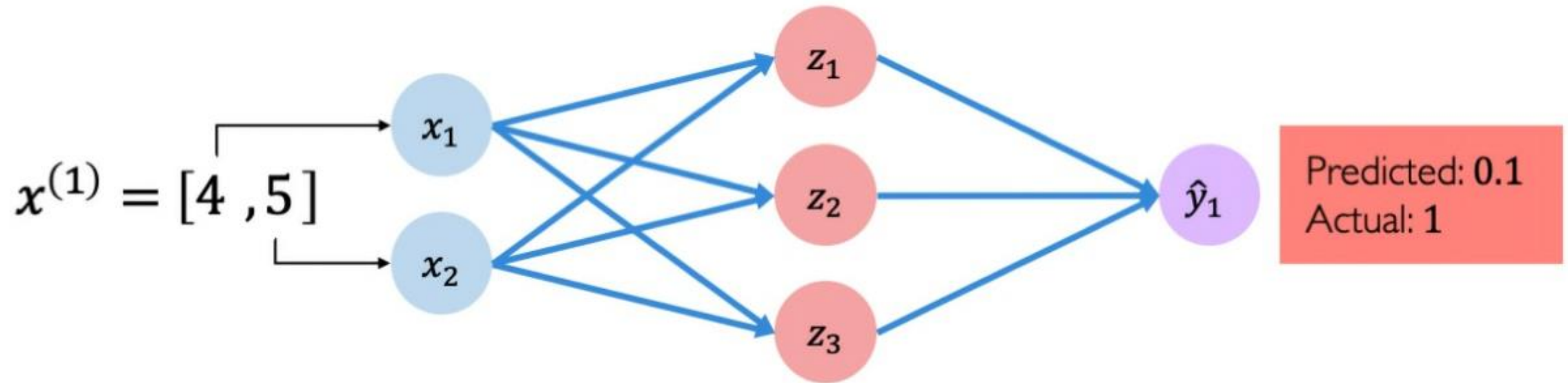


## Example Problem: Will I Pass this Class?



# Quantifying Loss

The Loss of our network measures the cost incurred from incorrect predictions



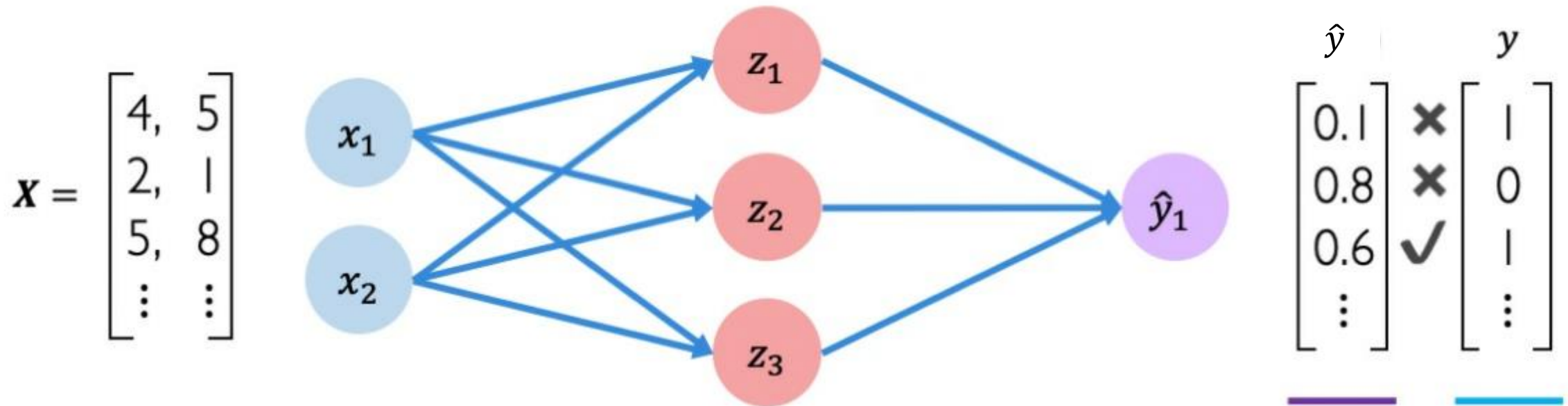
$$\mathcal{L}(\hat{y}, y)$$

Predicted    Actual



# Empirical Loss

The **empirical loss** measures the total loss over entire dataset



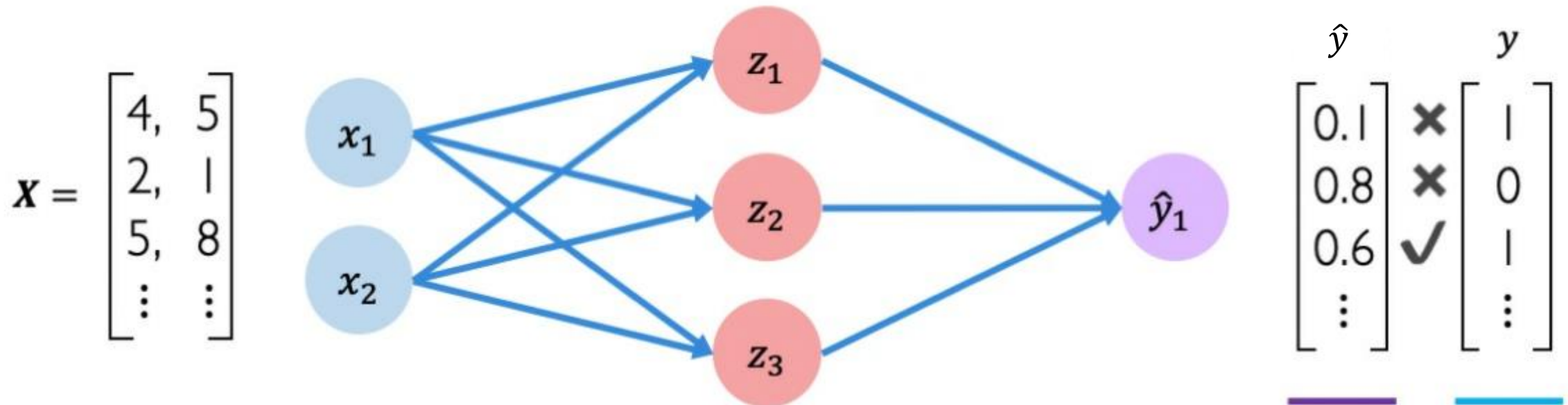
Also known as:

- Objective function
- Cost function
- Empirical Risk

$$J(W, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

# Binary Cross Entropy Loss

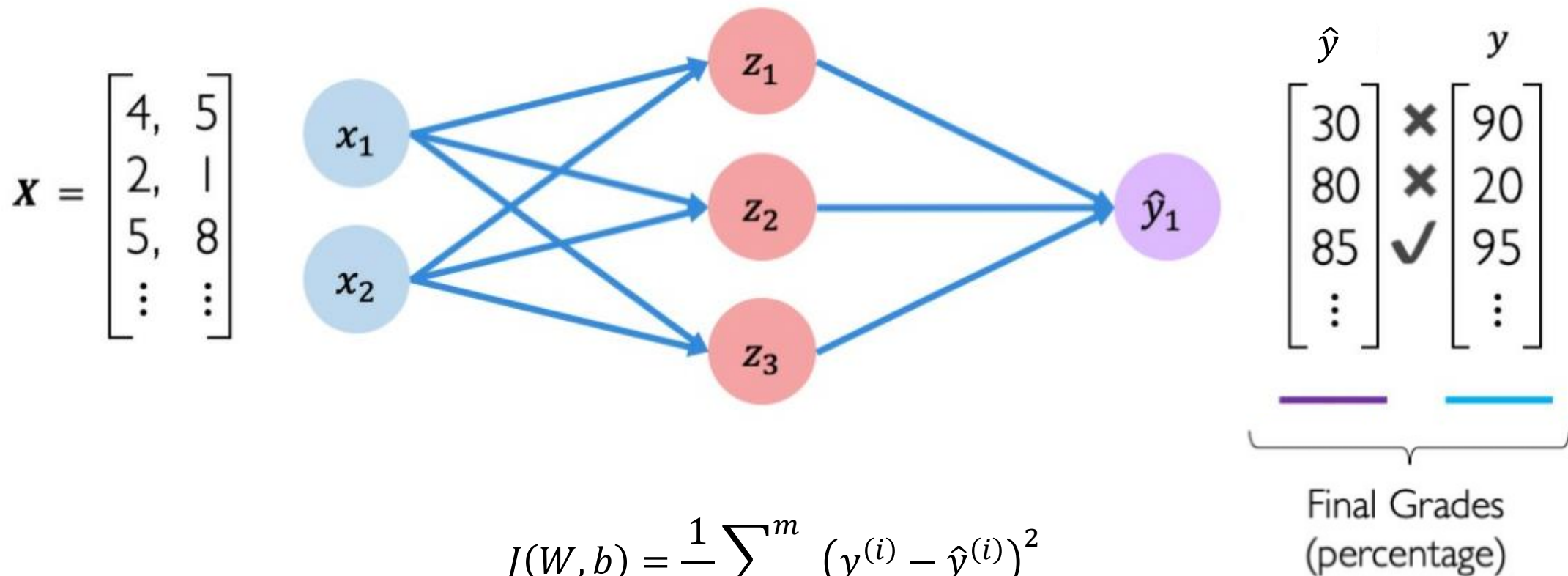
**Binary Cross entropy loss** can be used with models that output a probability between 0 and 1.



$$J(W, b) = \frac{1}{m} \sum_{i=1}^m -(y^{(i)} \ln(\hat{y}^{(i)}) + (1 - y^{(i)}) \ln(1 - \hat{y}^{(i)}))$$

# Mean Squared Error Loss

**Mean squared error loss** can be used with regression models that output continuous real numbers.



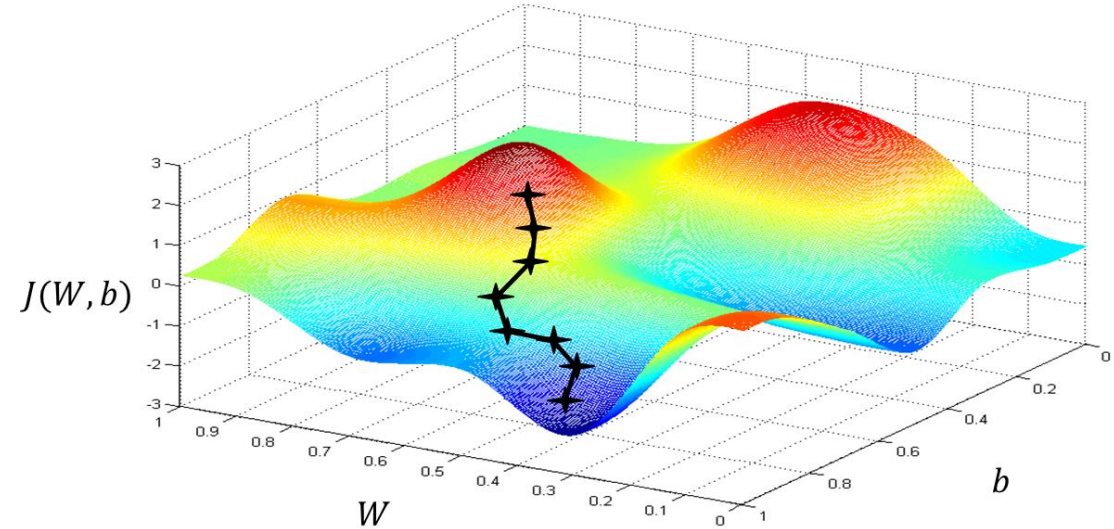
# Training Neural Networks

# Loss Optimization

We want to find the network weights and biases that achieve the lowest loss

$$J(W, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

$$[w^*, b^*] = \min_{W, b} J(W, b)$$

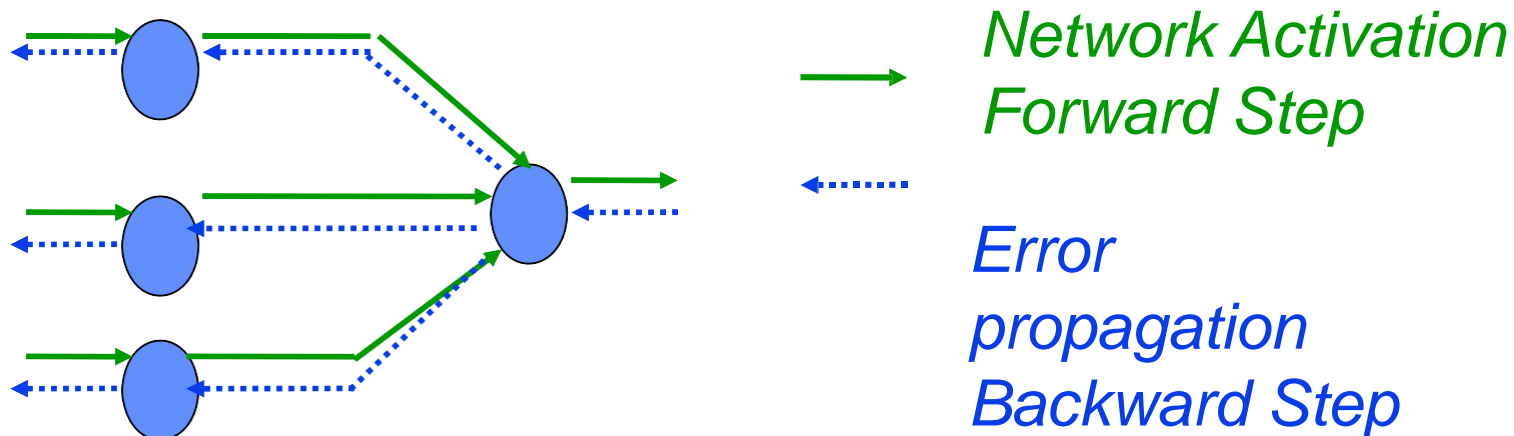


# Training a Neural Networks (computing weights)

## Random initialization of the weights

Backpropagation consists of the repeated application of the following two steps:

- Forward step: in this step the network is activated on one example and the error of (each neuron of) the output layer is computed.
- Backward step: in this step the network error is used for updating the weights. Starting at the output layer, the error is propagated backwards through the network, layer by layer. This is done by recursively computing the local gradient of each neuron.



# Derivatives of Activation Functions

## Derivatives

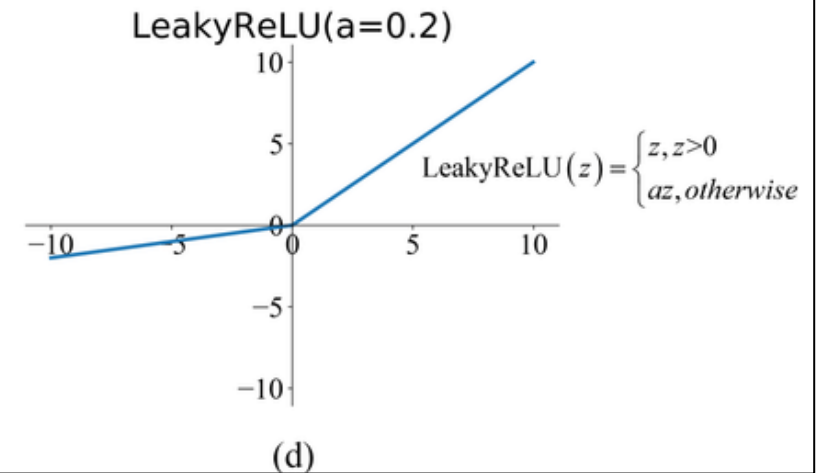
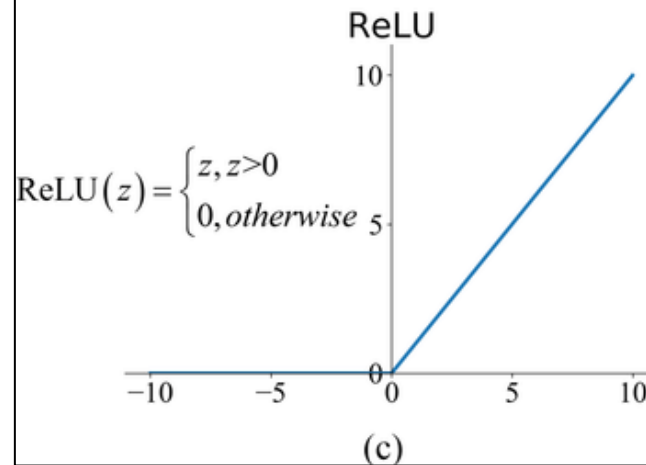
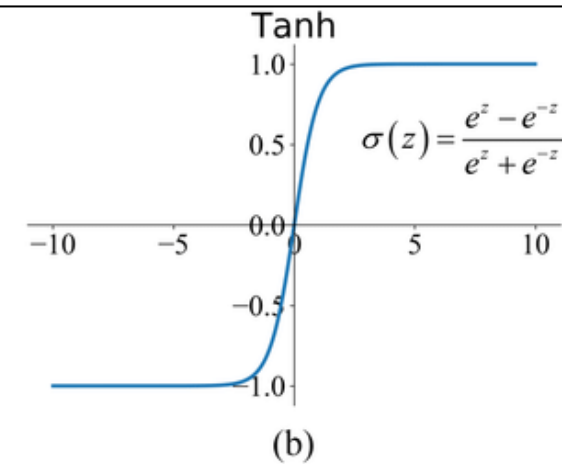
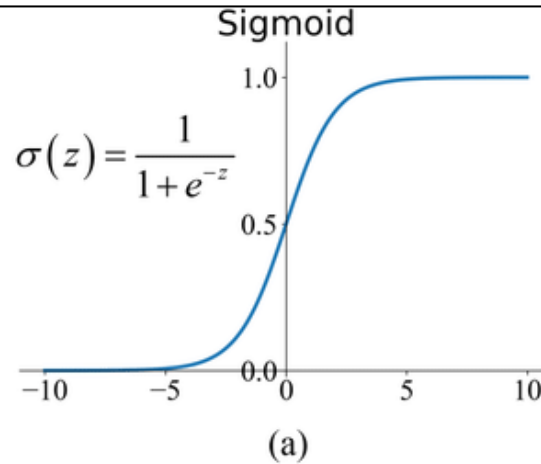
$$a) \quad \frac{d \sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$$

$$b) \quad \frac{d \sigma(x)}{dx} = 1 - \sigma(x)^2$$

$$c) \quad \frac{d \sigma(x)}{dx} = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \\ \text{undefine} & \text{if } z = 0^* \end{cases}$$

$$d) \quad \frac{d \sigma(x)}{dx} = \begin{cases} a & \text{if } z < 0 \\ 1 & \text{if } z > 0 \\ \text{undefine} & \text{if } z = 0^* \end{cases}$$

\* No relevant in practical scenarios



# Matrix – Partial Derivatives

---

The derivative of a vector function (a vector whose components are functions)  $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_m]^T$ , with respect to an input vector,  $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^T$  is written as

$$\frac{d\mathbf{y}}{d\mathbf{x}} = \begin{bmatrix} \frac{dy_1}{dx_1} & \dots & \frac{dy_1}{dx_n} \\ \vdots & \ddots & \vdots \\ \frac{dy_m}{dx_1} & \dots & \frac{dy_m}{dx_n} \end{bmatrix}$$

The derivative of a vector function  $\mathbf{y}$  with respect to an input vector,  $\mathbf{x}$  whose components represent a space is known as the Jacobian Matrix.



# Computing Gradients

## Logistic Regression

$$z = w^T x + b$$

$$a = \sigma(z)$$

$$\text{Cost}(a, y) = \mathcal{L}(a, y) = -(y \ln(a) + (1 - y) \ln(1 - a))$$

Derivatives:

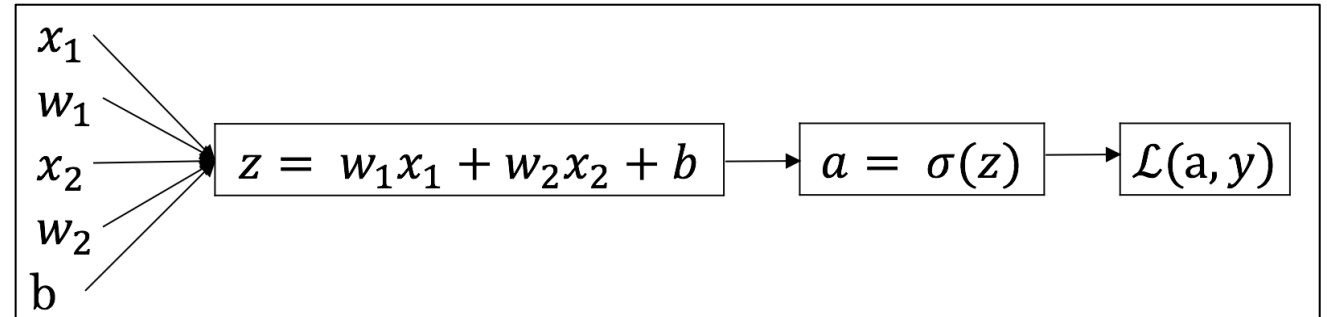
$$\frac{d\mathcal{L}}{da} = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\frac{d\mathcal{L}}{dz} = \frac{d\mathcal{L}}{da} \frac{da}{dz} = \left( -\frac{y}{a} + \frac{1-y}{1-a} \right) (a(1-a)) = a - y$$

$$\frac{d\mathcal{L}}{dw_1} = \frac{d\mathcal{L}}{da} \frac{da}{dz} \frac{dz}{dw_1} = (a - y)x_1$$

$$\frac{d\mathcal{L}}{dw_2} = \frac{d\mathcal{L}}{da} \frac{da}{dz} \frac{dz}{dw_2} = (a - y)x_2$$

$$\frac{d\mathcal{L}}{db} = \frac{d\mathcal{L}}{da} \frac{da}{dz} \frac{dz}{db} = (a - y)$$



Gradient Descent Algorithm for Logistic Regression:

$$w_1 := w_1 - \alpha \frac{d J(w, b)}{dw_1}$$

$$w_2 := w_2 - \alpha \frac{d J(w, b)}{dw_2}$$

$$b := b - \alpha \frac{d J(w, b)}{db}$$

# Neural Network Gradients

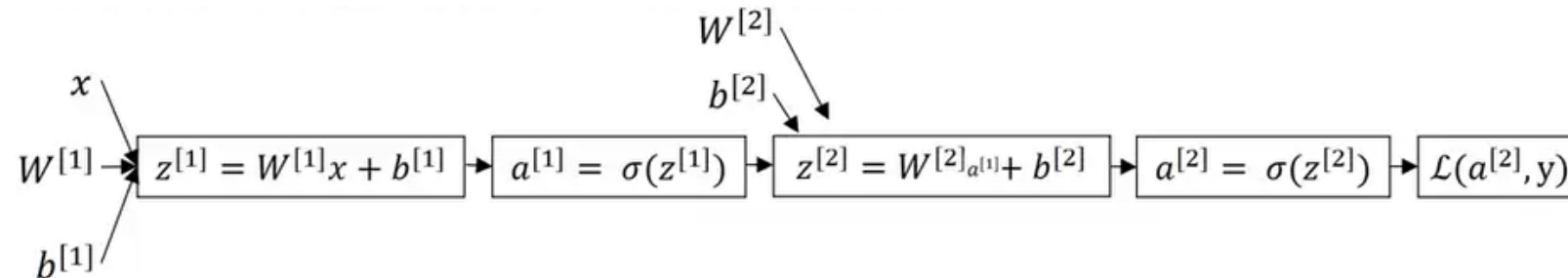
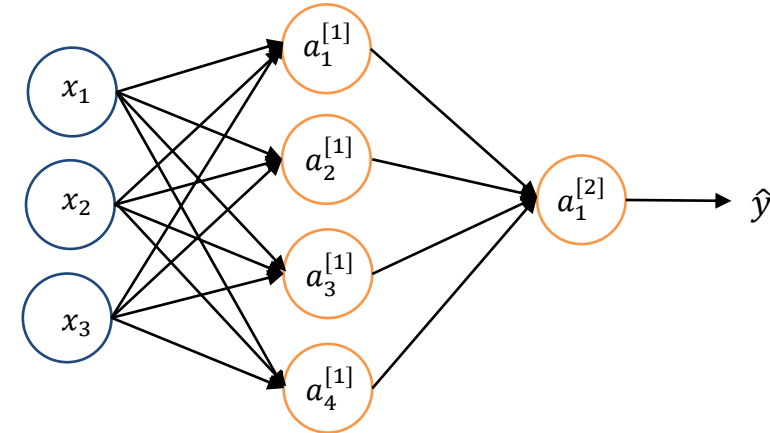
Let's suppose to design a Neural Network for Binary Classification

Parameters:  $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$

Cost Function:  $J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$

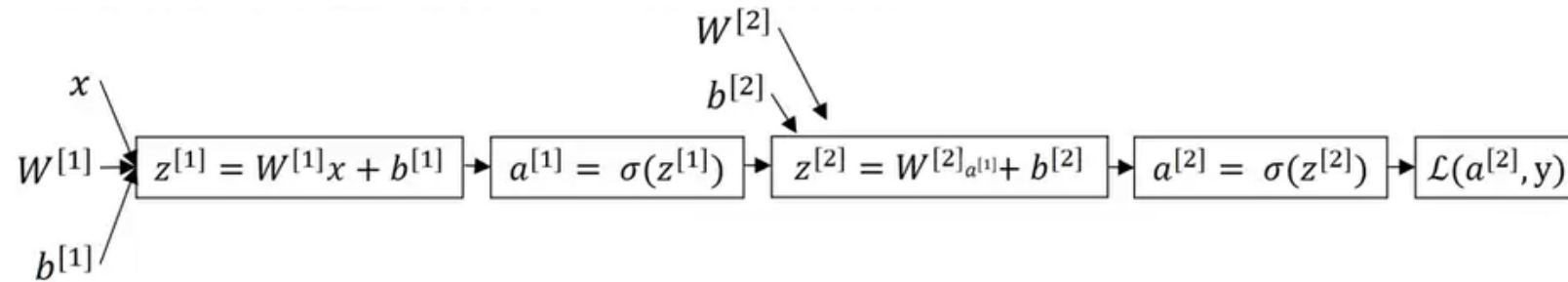
$\mathcal{L}(\hat{y}, y) = \mathcal{L}(a^{[2]}, y) = -(y \ln(a^{[2]}) + (1 - y) \ln(1 - a^{[2]}))$

Let's start with one sample:



Remember: Learning Algorithm - Forward and Back propagation steps.

# Neural Network Gradients



Forward propagation step:

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = g^{[2]}(z^{[2]})$$

# Neural Network Gradients

Loss:

$$\mathcal{L}(\hat{y}, y) = \mathcal{L}(a^{[2]}, y) = -(y \ln(a^{[2]}) + (1 - y) \ln(1 - a^{[2]}))$$

Back propagation step:

$$\frac{d\mathcal{L}}{da^{[2]}} = -\frac{y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}}$$

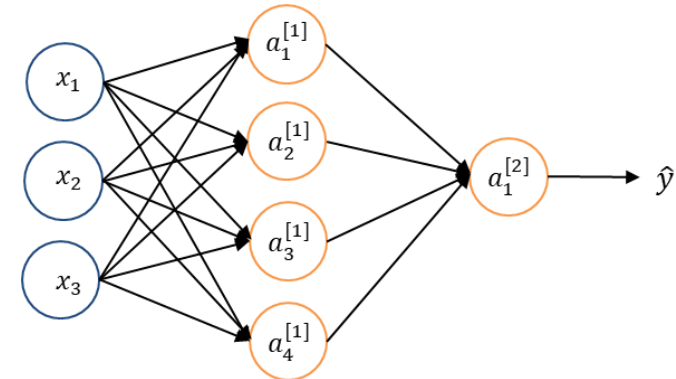
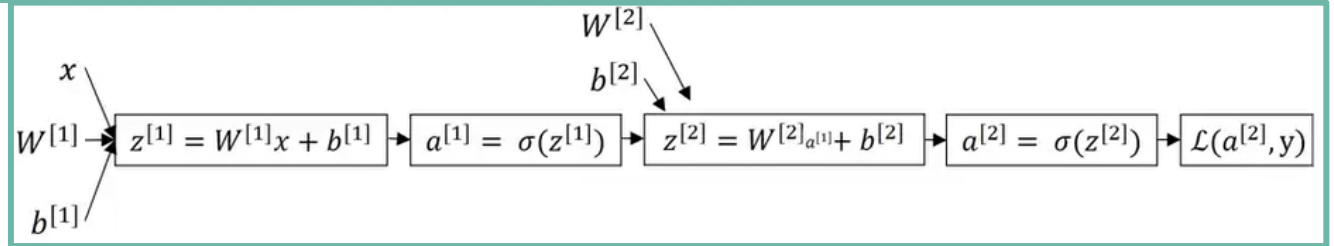
$$\frac{d\mathcal{L}}{dz^{[2]}} = \frac{d\mathcal{L}}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} = \left( -\frac{y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}} \right) (a^{[2]}(1-a^{[2]})) = a^{[2]} - y$$

$$\frac{d\mathcal{L}}{dW^{[2]}} = \frac{d\mathcal{L}}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} \frac{dz^{[2]}}{dW^{[2]}} = (a^{[2]} - y) a^{[1]T}$$

$$\frac{d\mathcal{L}}{db^{[2]}} = \frac{d\mathcal{L}}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} \frac{dz^{[2]}}{db^{[2]}} = (a^{[2]} - y)$$

$$\frac{d\mathcal{L}}{da^{[1]}} = \frac{d\mathcal{L}}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} \frac{dz^{[2]}}{da^{[1]}} = (a^{[2]} - y) W^{[2]T} = W^{[2]T} (a^{[2]} - y)$$

$$\frac{d\mathcal{L}}{dz^{[1]}} = \frac{d\mathcal{L}}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} \frac{dz^{[2]}}{da^{[1]}} \frac{da^{[1]}}{dz^{[1]}} = W^{[2]T} (a^{[2]} - y) * (a^{[1]} * (1 - a^{[1]}))$$



**Given in input X:**

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

(4,1) (4,3) (3,1) (4,1) - Dimensionality

$$a^{[1]} = g^{[1]}(z^{[1]})$$

(4,1) (4,1) - Dimensionality

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

(1,1) (1,4) (4,1) (1,1) - Dimensionality

$$a^{[2]} = g^{[2]}(z^{[2]})$$

(1,1) (1,1) - Dimensionality

# Neural Network Gradients

Back propagation step:

$$\frac{d\mathcal{L}}{da^{[2]}} = -\frac{y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}}$$

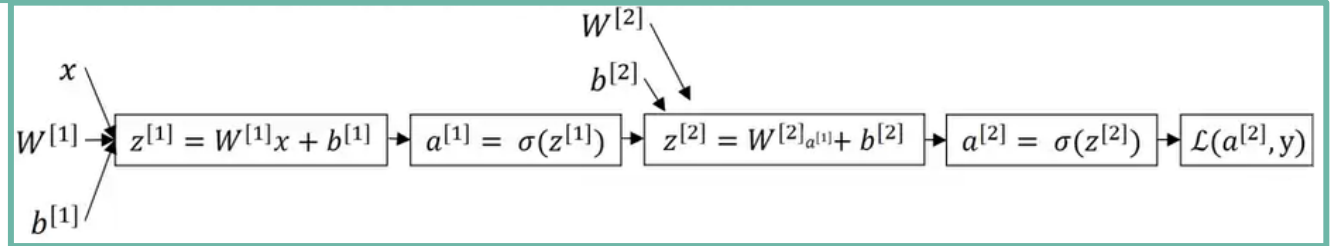
$$\frac{d\mathcal{L}}{dz^{[2]}} = \frac{d\mathcal{L}}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} = \left( -\frac{y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}} \right) (a^{[2]}(1-a^{[2]})) = a^{[2]} - y$$

$$\frac{d\mathcal{L}}{dW^{[2]}} = \frac{d\mathcal{L}}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} \frac{dz^{[2]}}{dW^{[2]}} = (a^{[2]} - y) a^{[1]T}$$

$$\frac{d\mathcal{L}}{db^{[2]}} = \frac{d\mathcal{L}}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} \frac{dz^{[2]}}{db^{[2]}} = (a^{[2]} - y)$$

$$\frac{d\mathcal{L}}{da^{[1]}} = \frac{d\mathcal{L}}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} \frac{dz^{[2]}}{da^{[1]}} = (a^{[2]} - y) W^{[2]T} = W^{[2]T} (a^{[2]} - y)$$

$$\frac{d\mathcal{L}}{dz^{[1]}} = \frac{d\mathcal{L}}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} \frac{da^{[1]}}{dz^{[1]}} = W^{[2]T} (a^{[2]} - y) * (a^{[1]} * (1 - a^{[1]}))$$



$$\frac{d\mathcal{L}}{dW^{[1]}} = \frac{d\mathcal{L}}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} \frac{dz^{[2]}}{da^{[1]}} \frac{da^{[1]}}{dz^{[1]}} \frac{dz^{[1]}}{dW^{[1]}} = \frac{d\mathcal{L}}{dz^{[1]}} x^T$$

$$\frac{d\mathcal{L}}{db^{[1]}} = \frac{d\mathcal{L}}{da^{[2]}} \frac{da^{[2]}}{dz^{[2]}} \frac{da^{[1]}}{dz^{[1]}} \frac{dz^{[1]}}{db^{[1]}} = \frac{d\mathcal{L}}{dz^{[1]}}$$

Gradient Descent Algorithm :

$$W^{[1]} := W^{[1]} - \alpha \frac{dJ(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]})}{dW^{[1]}}$$

...

# Neural Network Gradients with $m$ examples

$$J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(a^{[2](i)}, y^{(i)})$$

Back propagation step with  $m$  examples:

$$\frac{dJ(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]})}{dW^{[2]}} = \frac{1}{m} \sum_{i=1}^m \frac{d\mathcal{L}(a^{[2](i)}, y^{(i)})}{da^{[2](i)}}$$

$$\frac{dJ(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]})}{db^{[2]}} = \frac{1}{m} \sum_{i=1}^m \frac{d\mathcal{L}(a^{[2](i)}, y^{(i)})}{da^{[2](i)}}$$

$$\frac{dJ(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]})}{dW^{[1]}} = \frac{1}{m} \sum_{i=1}^m \frac{d\mathcal{L}(a^{[2](i)}, y^{(i)})}{da^{[1](i)}}$$

$$\frac{dJ(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]})}{db^{[1]}} = \frac{1}{m} \sum_{i=1}^m \frac{d\mathcal{L}(a^{[2](i)}, y^{(i)})}{da^{[1](i)}}$$

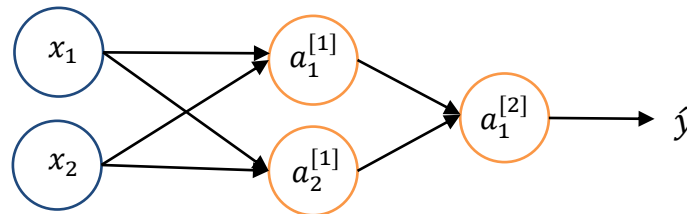
Gradient Descent Algorithm:

$$\begin{aligned} W^{[2]} &:= W^{[2]} - \alpha \frac{dJ(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]})}{dW^{[2]}} \\ b^{[2]} &:= b^{[2]} - \alpha \frac{dJ(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]})}{db^{[2]}} \\ &\dots \end{aligned}$$

# Random Initialization

What happens if you initialize weights to zero?

Let's start with this example:



**Derivatives from previous slides:**

$$\frac{d\mathcal{L}}{dW^{[2]}} = (a^{[2]} - y)a^{[1]T}$$

$$\frac{d\mathcal{L}}{dz^{[1]}} = W^{[2]T}(a^{[2]} - y) * (a^{[1]} * (1 - a^{[1]}))$$

$$\frac{d\mathcal{L}}{dW^{[1]}} = \frac{d\mathcal{L}}{dz^{[1]}} x^T$$

If  $W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$ ,  $b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$  and,  $W^{[2]} = \begin{bmatrix} 0 & 0 \end{bmatrix}$  for any input sample in forward propagation  $a_1^{[1]} = a_2^{[1]} \rightarrow \frac{d\mathcal{L}}{dW_1^{[2]}} = \frac{d\mathcal{L}}{dW_2^{[2]}} \rightarrow \frac{d\mathcal{L}}{dz_1^{[1]}} = \frac{d\mathcal{L}}{dz_2^{[1]}}$

Let's suppose  $\frac{d\mathcal{L}}{dz_1^{[1]}} = \frac{d\mathcal{L}}{dz_2^{[1]}} = v$ , when we compute the  $\frac{d\mathcal{L}}{dW^{[1]}}$  we can see that:

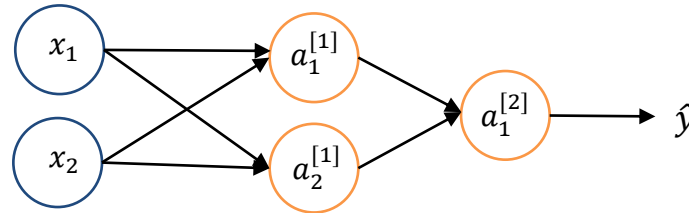
$$\frac{d\mathcal{L}}{dW^{[1]}} = \frac{d\mathcal{L}}{dz^{[1]}} x^T = \begin{bmatrix} v \\ v \end{bmatrix} \begin{bmatrix} x_1 & x_2 \end{bmatrix} = \begin{bmatrix} vx_1 & vx_2 \\ vx_1 & vx_2 \end{bmatrix}$$

Thus, in the second forward interaction computing we will have  $a_1^{[1]} = a_2^{[1]}$  and so on.... It's called "symmetry problem"

Random initialization is a solution! (Note: biases initialization to zero is not a problem)

# Random Initialization

Let's start with this example:



$$W^{[1]} = np.random.randn((2,2)) * \gamma$$

$$b^{[1]} = np.zeros((2,1))$$

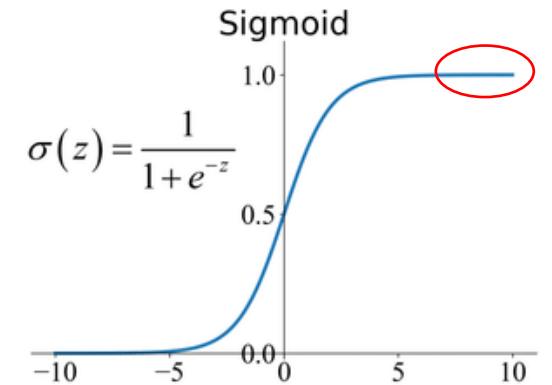
$$W^{[2]} = np.random.randn((1,2)) * \gamma$$

$$b^{[2]} = 0$$

For example,  $\gamma = 0.01$

The main idea is to try to keep values small!

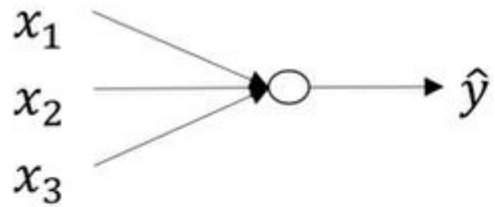
Otherwise, if you use **large weights**, for example with a **sigmoid function** in the last layer, your output values will go on the flat region of the function where **the slop or the gradient is very small. Gradient Descent will be very slow.**



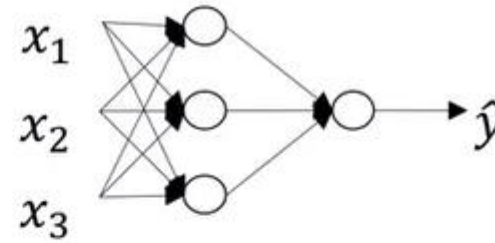


# What is a Deep Neural Network?

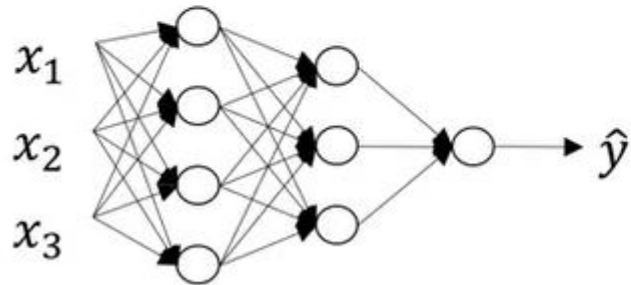
«Shallow Network»



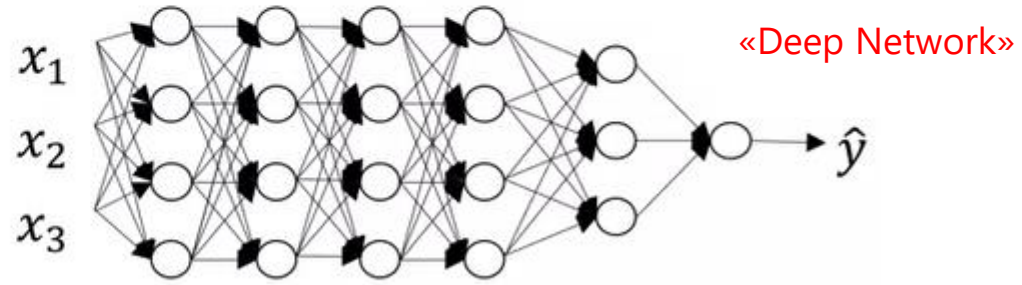
logistic regression



1 hidden layer

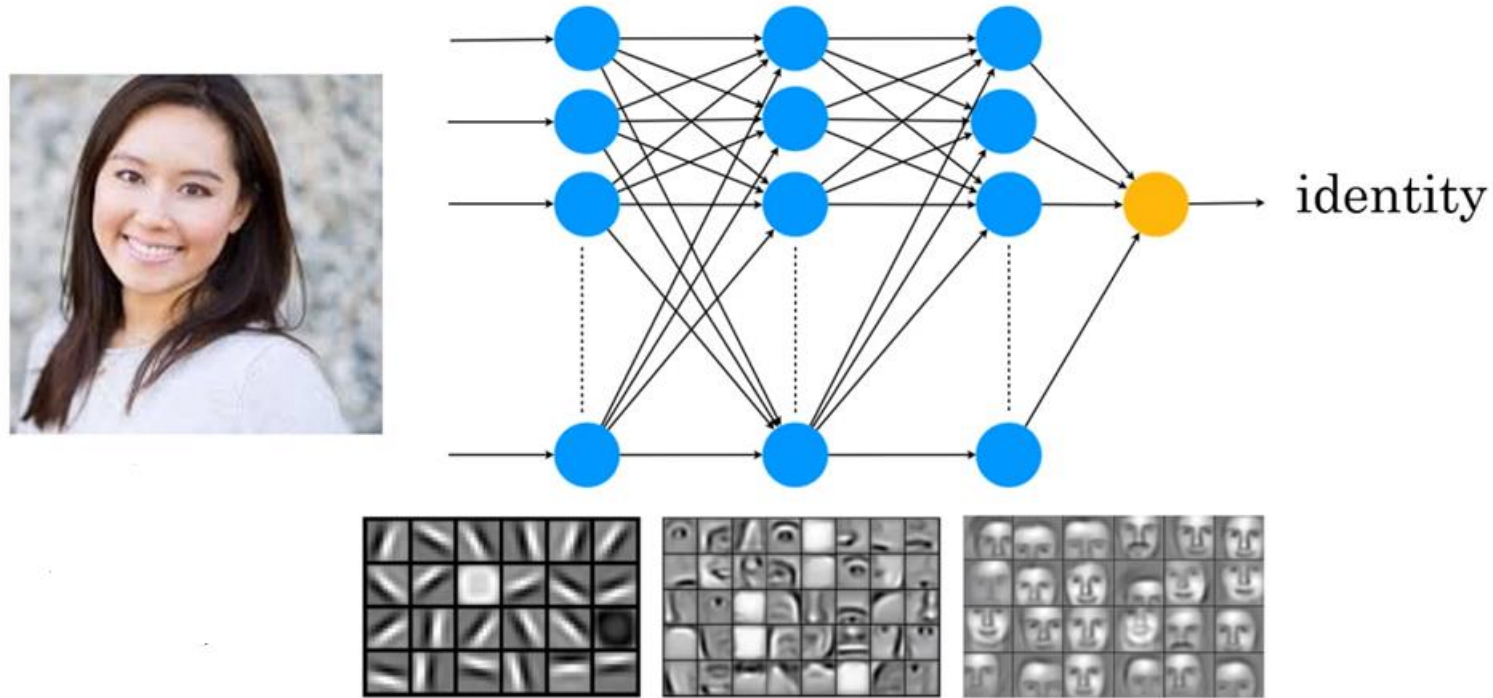


2 hidden layers



5 hidden layers

# Intuition about deep representation



We will better understand this image when we are using Convolution Neural Networks.

The idea could be applied to audio context: phonemes... words... sentences

# What are hyperparameters?

Parameters:  $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots$

Hyperparameters:

- Learning rate
- Number of interaction (epochs)
- Number of hidden layers
- Number of hidden units in each hidden layers
- Activation functions

Right now, this process is empirically.

