# Neural Networks

Dr. Amani RAAD

amani.raad@ul.edu.lb
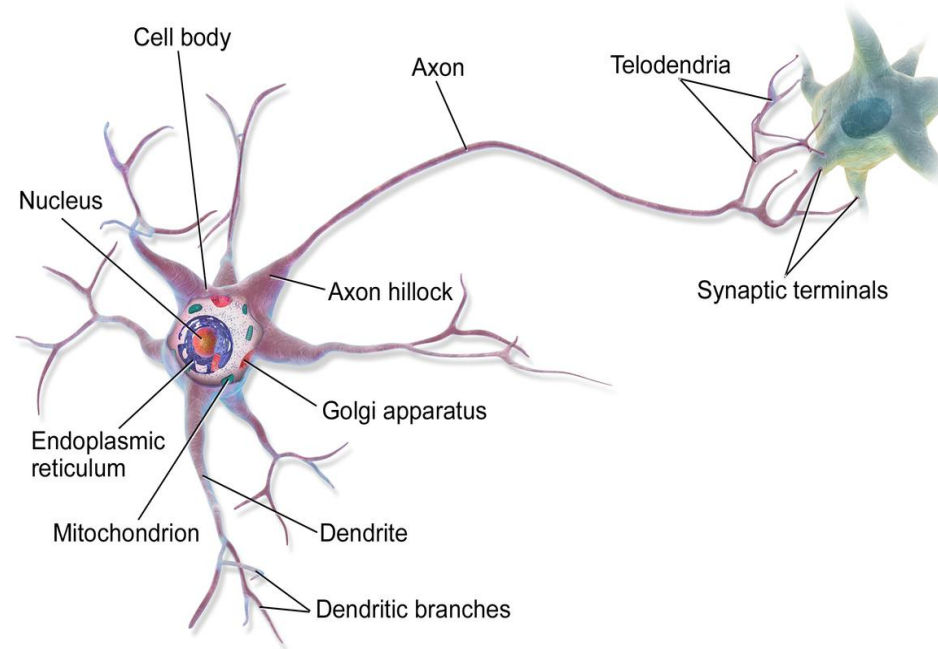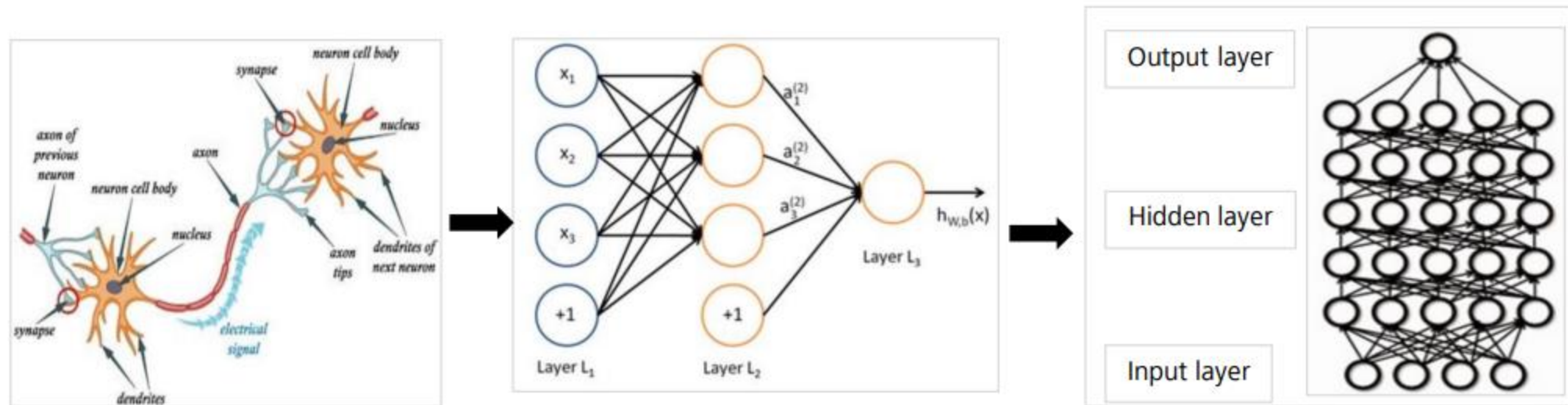
2022-2023

# Neural Network Definition

- Currently, the definition of the neural network is not determined yet. Hecht Nielsen, a neural network researcher in the US, defines a neural network as "a computing system made up of a number of simple, highly interconnected processing elements, which process information by their dynamic state response to external inputs."

- Based on the origin, features, and interpretations of the neural network, it can be simply defined as **an information processing system designed to simulate human brain's structure and functions**.

- **Artificial neural network (neural network for short):** refers to a network composed of artificial neurons. It abstracts and simplifies a human brain based on its microscopic structure and functions. It is an important way to simulate human intelligence and reflects some basic features of human brain functions, such as parallel information processing, learning, association, pattern classification, and memory.

# Biologic Neuron

- The basic computational unit of the brain is a neuron

- Neurons receive input signal from dendrites and produce output signal along axon, which interact with the dendrites of other neurons via synaptic weights

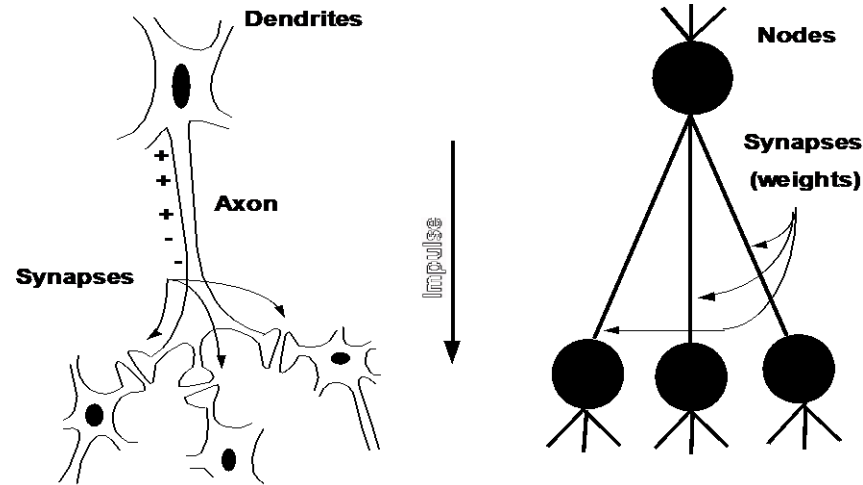- The neuron collects info through dendrites; if this is large enough the neuron fires.

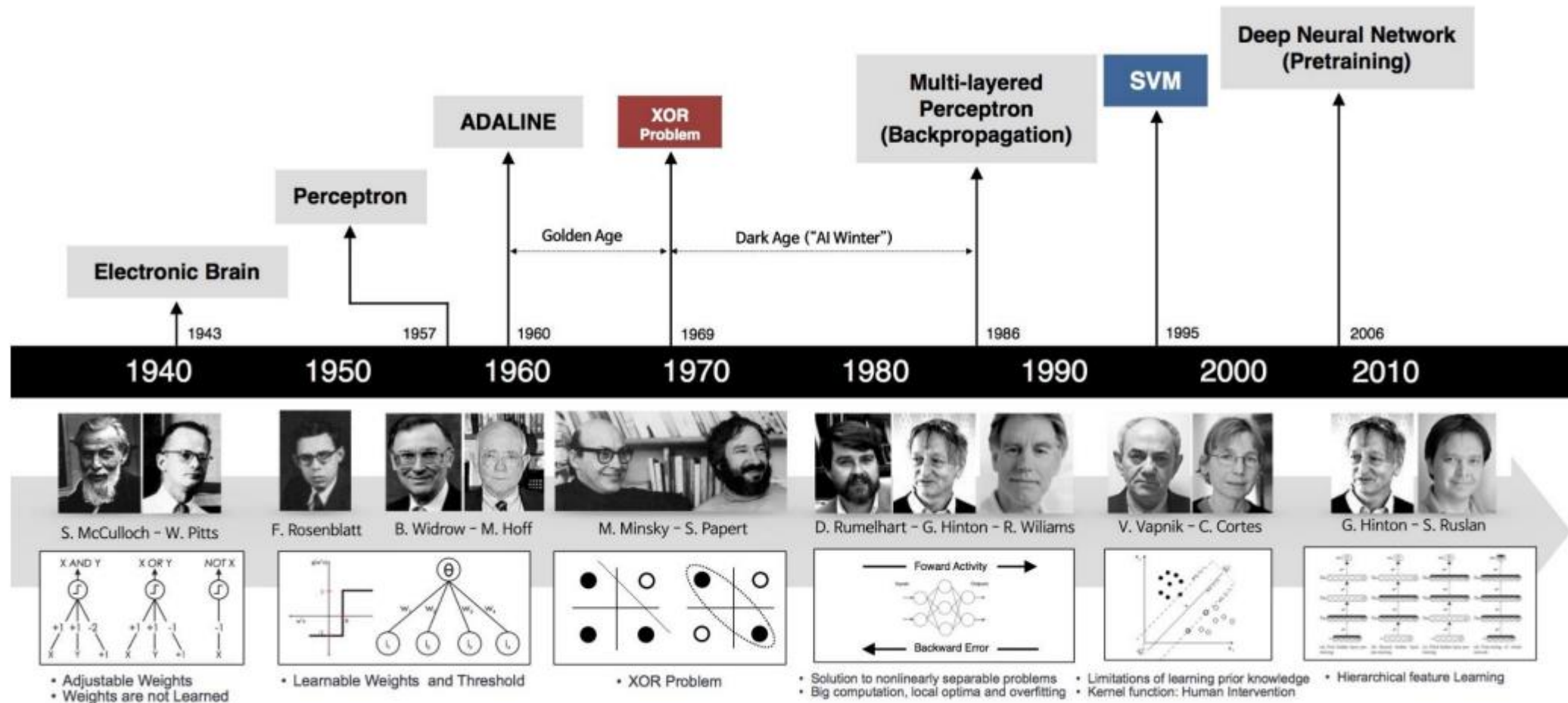# From the brain to the perceptron to the Multi Layer Perceptrons

# Connectionist Models

- Consider humans:
  - Neuron switching time
    ~ 0.001 second
  - Number of neurons
    ~ $10^{10}$
  - Connections per neuron
    ~ $10^{4\text{-}5}$
  - Scene recognition time
    ~ 0.1 second
  - 100 inference steps doesn't seem like enough
    $\rightarrow$ much parallel computation

- Properties of artificial neural nets (ANN)
  - Many neuron-like threshold switching units
  - Many weighted interconnections among units
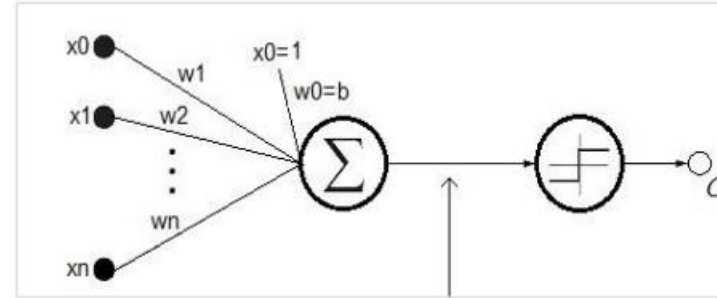  - Highly parallel, distributed processes

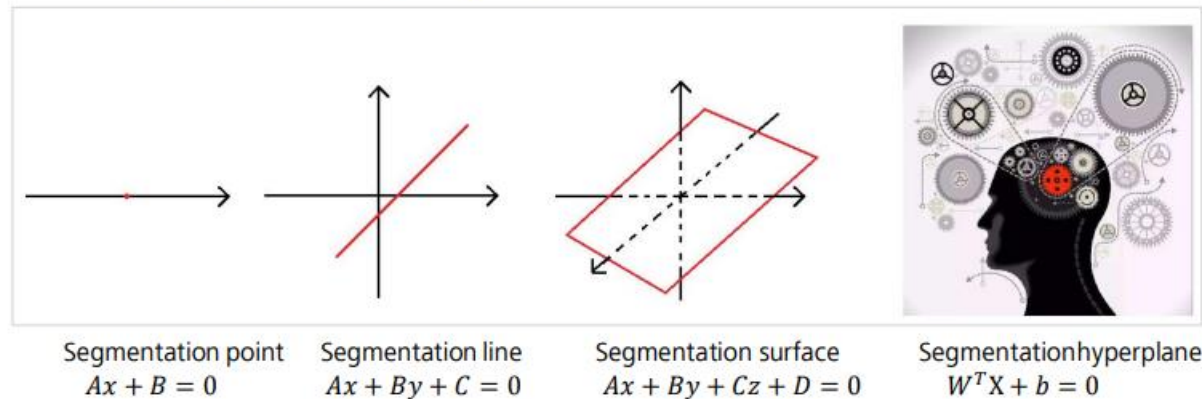# Neural Networks Milestones (Deep Learning)

# Perceptron



- **Input vector:** $X = [x_0, x_1, \ldots, x_n]^T$.

- **Weight:** $W = [\omega_0, \omega_1, \ldots, \omega_n]^T$, where $\omega_0$ is the bias.

- **Activation function:** $O = sign(net) = \begin{cases} 1, net > 0, \\ -1, otherwise. \end{cases}$

$$net = \sum_{i=0}^{n} \omega_i x_i = W^T X$$

- The perceptron is equivalent to a classifier. Its input is the high-dimensional vector $X$ and it performs binary classification on input samples in the high-dimensional space. If $W^T X > 0$, $o = 1$ and the sample is classified into one class. Otherwise, $o = -1$ and the sample is classified into the other class. What is the boundary? $W^T X = 0$. This is a high-dimensional hyperplane.
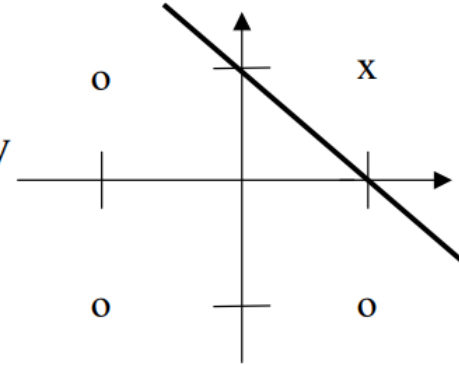


| Segmentation point | Segmentation line | Segmentation surface | Segmentationhyperplane |
|---|---|---|---|
| $Ax + B = 0$ | $Ax + By + C = 0$ | $Ax + By + Cz + D = 0$ | $W^T X + b = 0$ |

- Examples of linearly separable classes

  - Logical **AND** function

    patterns (bipolar) decision boundary

    | x1 | x2 | y | |
    |----|----|----|----|
    | -1 | -1 | -1 | w1 = 1 |
    | -1 | 1 | -1 | w2 = 1 |
    | 1 | -1 | -1 | b = -1 |
    | 1 | 1 | 1 | $\theta = 0$ |

    $-1 + x1 + x2 = 0$

    x: class I (y = 1)
    o: class II (y = -1)

  - Logical **OR** function

    patterns (bipolar) decision boundary

    | x1 | x2 | y | |
    |----|----|----|----|
    | -1 | -1 | -1 | w1 = 1 |
    | -1 | 1 | 1 | w2 = 1 |
    | 1 | -1 | 1 | b = 1 |
    | 1 | 1 | 1 | $\theta = 0$ |

    $1 + x1 + x2 = 0$

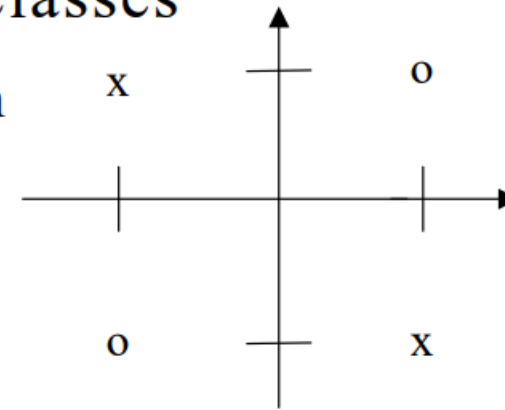    x: class I (y = 1)
    o: class II (y = -1)

- Examples of linearly inseparable classes
  - Logical **XOR** (exclusive OR) function

    patterns  (bipolar)  decision boundary

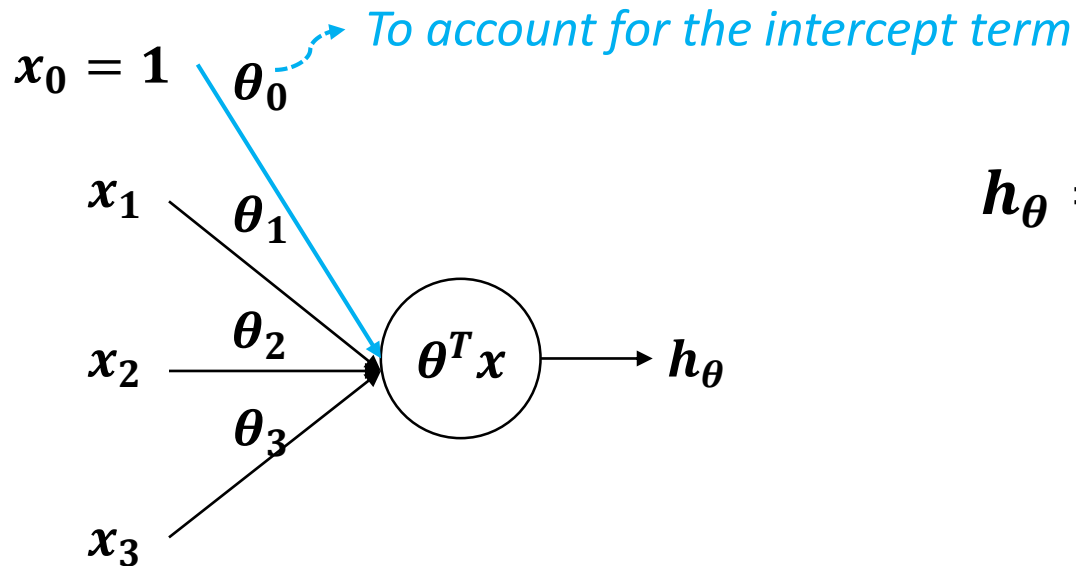    | x1 | x2 | y |
    |----|----|----|
    | -1 | -1 | -1 |
    | -1 | 1 | 1 |
    | 1 | -1 | 1 |
    | 1 | 1 | -1 |

  x: class I (y = 1)
  o: class II (y = -1)

  No line can separate these two classes, as can be seen from

# Linear Regression

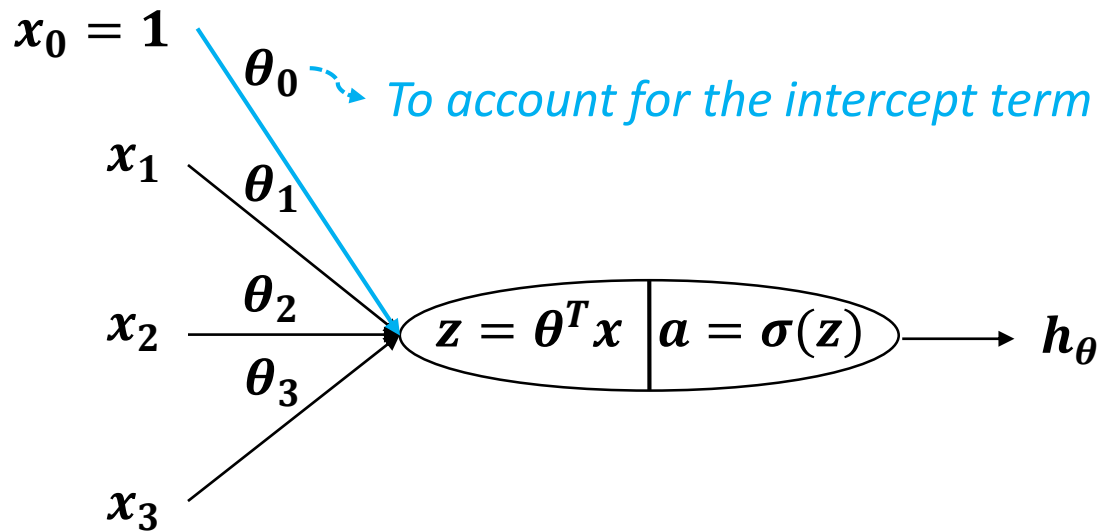- What is the hypothesis function of **linear regression**?
  - $h_\theta = \theta^T x$, where $\theta = [\theta_0, \theta_1, ..., \theta_m]$, $x = [x_0, x_1, ..., x_m]$, and $x_0 = 1$

*To account for the intercept term*

$x_0 = 1$   $\theta_0$

$x_1$   $\theta_1$

$x_2$   $\theta_2$   $\theta^T x$ → $h_\theta$

$\theta_3$

$x_3$

$$h_\theta = \theta^T x = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} [x_0 \; x_1 \; x_2 \; x_3]$$

$$= \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

# Logistic Regression

- What is the hypothesis function of **logistic regression**?
  - $h_\theta = \sigma(\theta^T x)$, where $\theta = [\theta_0, \theta_1, \ldots, \theta_m]$, $x = [x_0, x_1, \ldots, x_m]$, $x_0 = 1$, and
    $\sigma(z) = \frac{1}{1+e^{-z}}$

$$h_\theta = a = \sigma(z) = \sigma(\theta^T x) = \sigma\left(\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} [x_0 \ x_1 \ x_2 \ x_3]\right)$$

$x_0 = 1$

$\theta_0$ — *To account for the intercept term*

$x_1$
$\theta_1$

$x_2$ — $\theta_2$

$\theta_3$

$x_3$

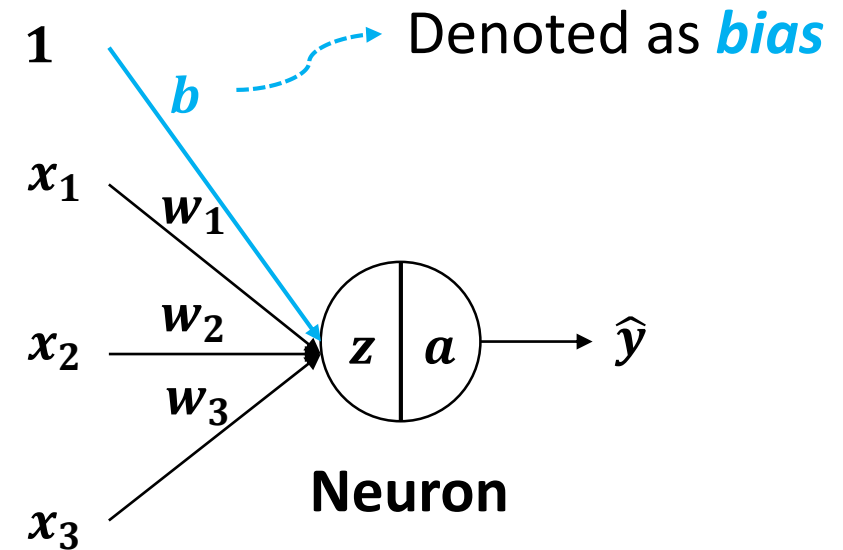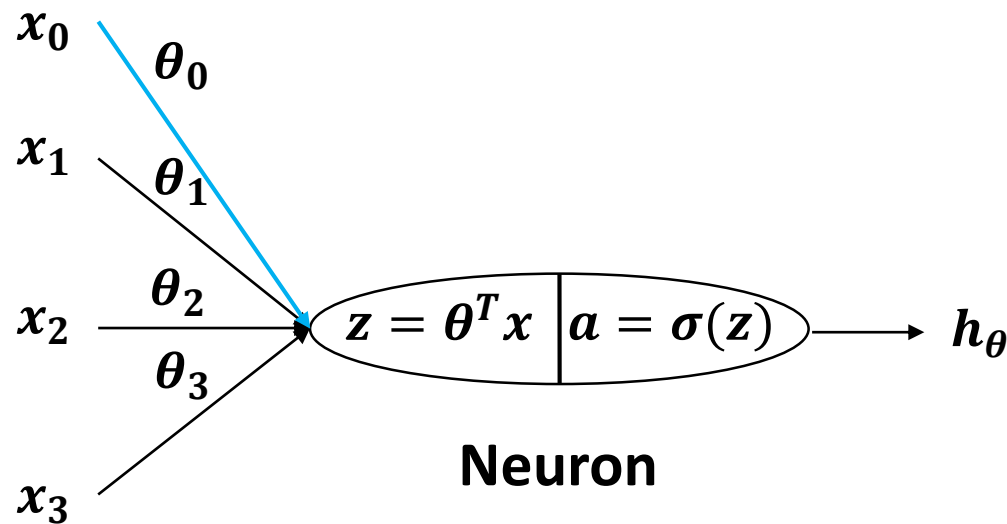$z = \theta^T x \mid a = \sigma(z)$ ⟶ $h_\theta$

$$= \sigma(\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3)$$

$$= \sigma(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3)$$

$$= \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3)}}$$
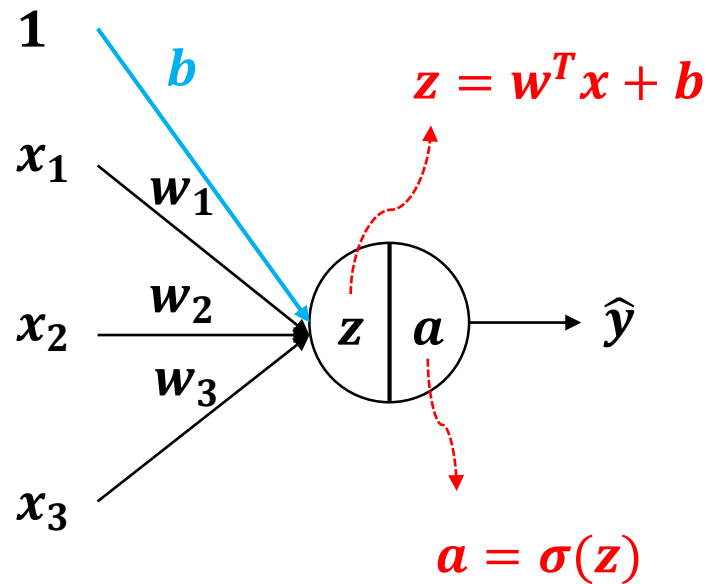
# Towards Neural Networks

- Technically, logistic regression is a **neural network** with only 1 neuron



Using the notations in the neural network literature, where $\boldsymbol{\theta} = \boldsymbol{w} = [\boldsymbol{w_1}, \boldsymbol{w_2}, \boldsymbol{w_3}]$ ($\boldsymbol{w_0}$ is not part of this vector here), $\boldsymbol{h_\theta} = \widehat{\boldsymbol{y}}$, and $\boldsymbol{\theta_0} = \boldsymbol{w_0} = \boldsymbol{b}$

# Towards Neural Networks

- Technically, logistic regression is a **neural network** with only 1 neuron

$$\hat{y} = a = \sigma(z) = \sigma(w^T x + b) = \sigma(\begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} [x_1 \; x_2 \; x_3] + b)$$

$$z = w^T x + b$$
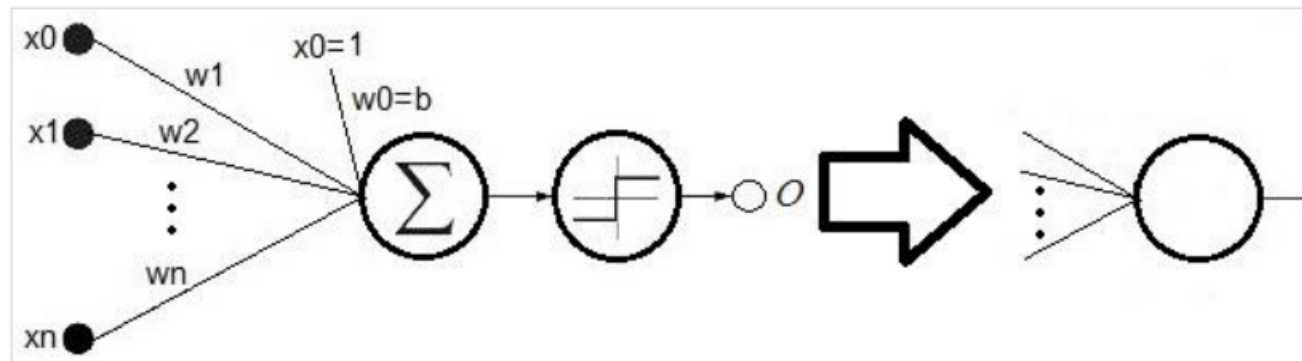
$$a = \sigma(z)$$

$$= \sigma(w_1 x_1 + w_2 x_2 + w_3 x_3 + b)$$

$$= \frac{1}{1 + e^{-(w_1 x_1 + w_2 x_2 + w_3 x_3 + b)}}$$
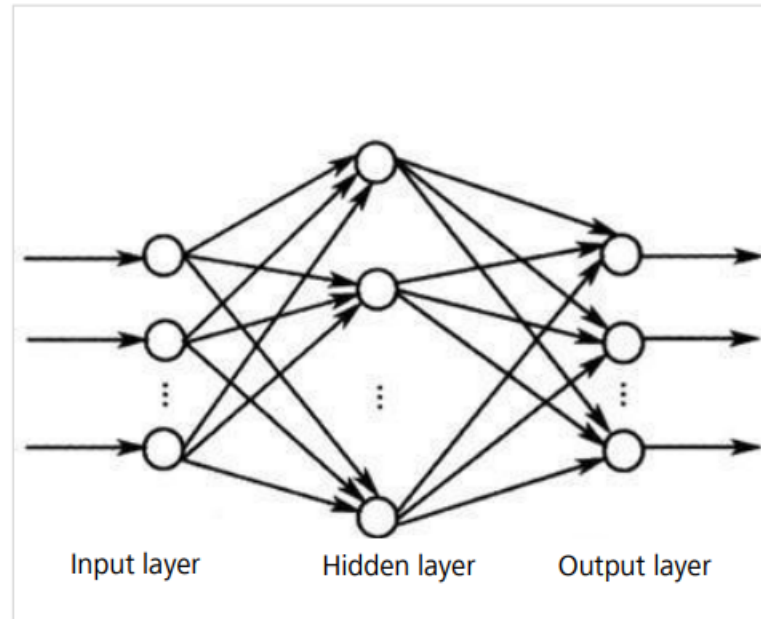
# Multi-layers perceptrons MLP
# ML Fully Connected

◆ A single perceptron has limited representation capability and can only represent the linear decision surface (hyperplane). If we connect many perceptrons like a human brain does and then replace the activation function with a non-linear function, we can express a wide range of non-linear surfaces.
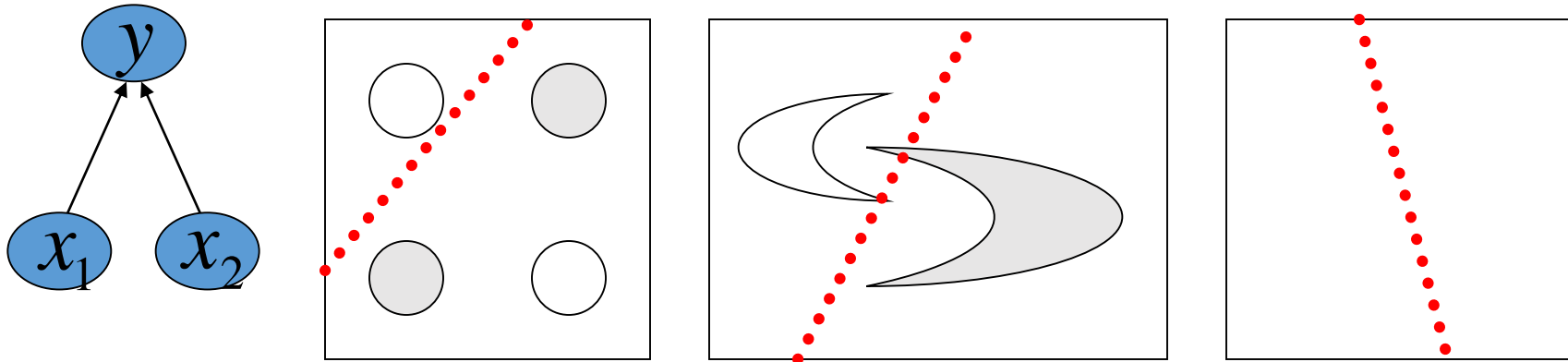
# Feed Forward Neural Network

- **The feedforward neural network** is one of the simplest neural networks in which neurons are arranged hierarchically. It is the most widely neural network with the fastest development.

- The input node does not support computation, but is merely used to represent each element value of the input vector.

- Each node represents a neuron that supports computation, which is called a computational unit. Each neuron is connected only to the neurons at the previous layer.

- A hidden layer receives the output from the previous layer and sends the results to the next layer. A unidirectional multi-layer structure is adopted. Each layer contains several neurons. The neurons at the same layer are not connected to each other. Inter-layer information is transmitted only in one direction.



Input layer      Hidden layer      Output layer
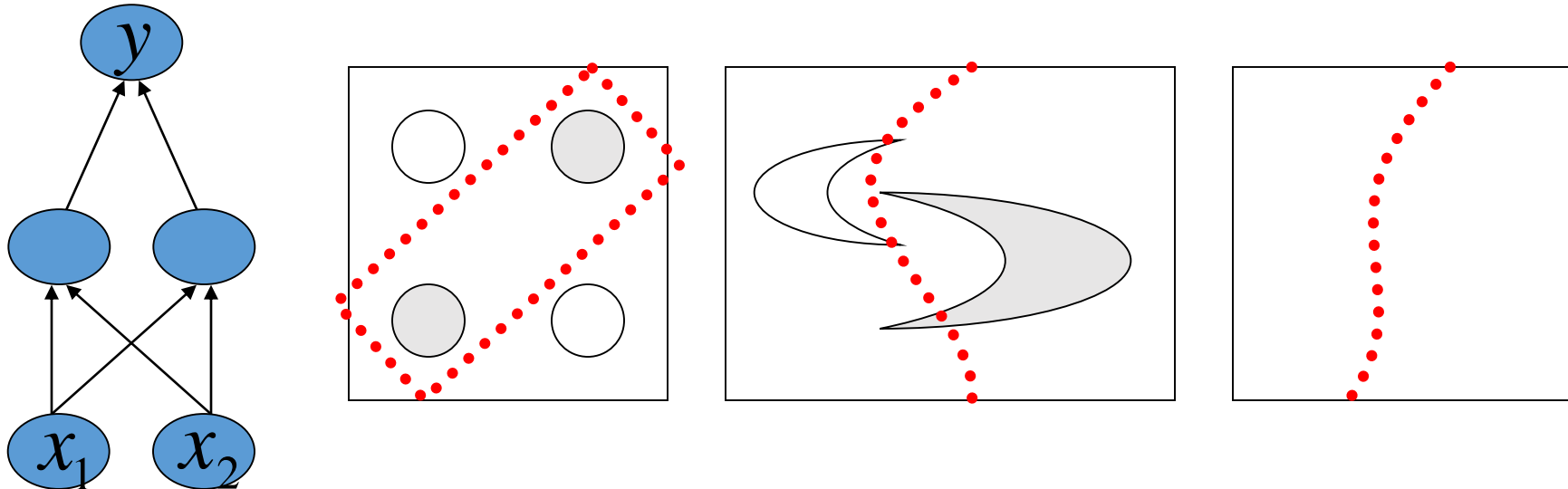
# Decision Boundary

- 0 hidden layers: linear classifier
  - Hyperplanes



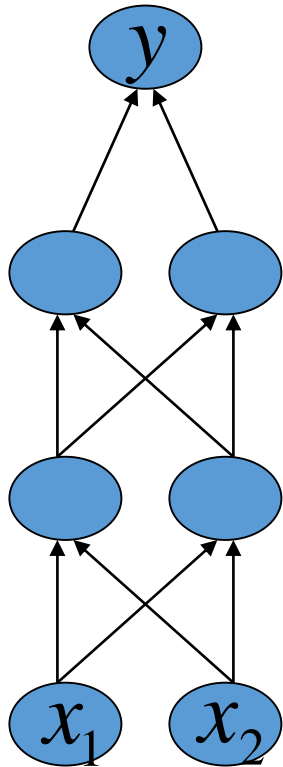Example from to Eric Postma via Jason Eisner

# Decision Boundary

- 1 hidden layer
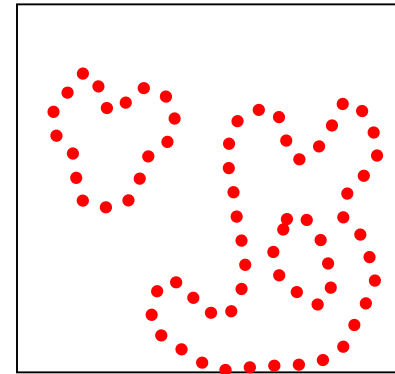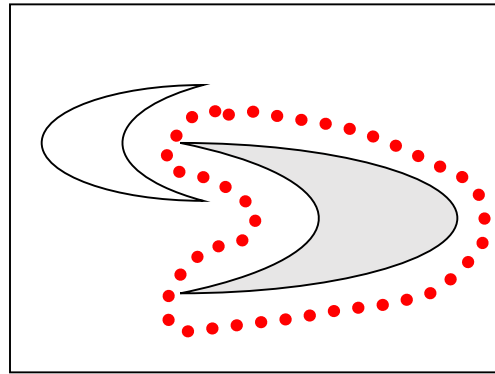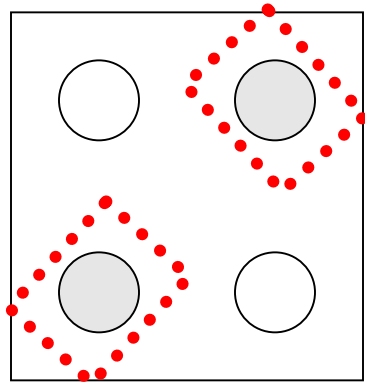    - Boundary of convex region (open or closed)



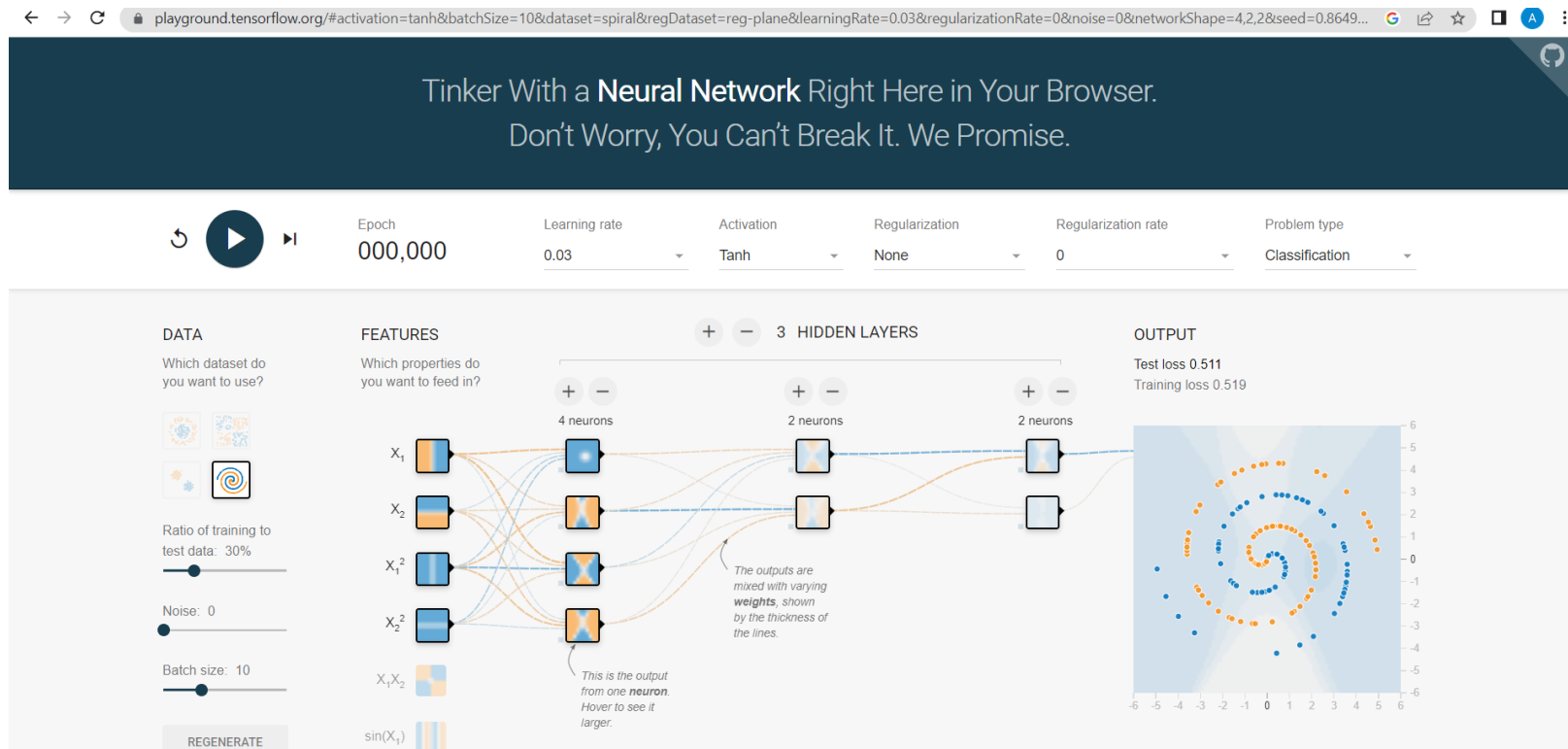Example from to Eric Postma via Jason Eisner

# Decision Boundary



- 2 hidden layers
  - Combinations of convex regions

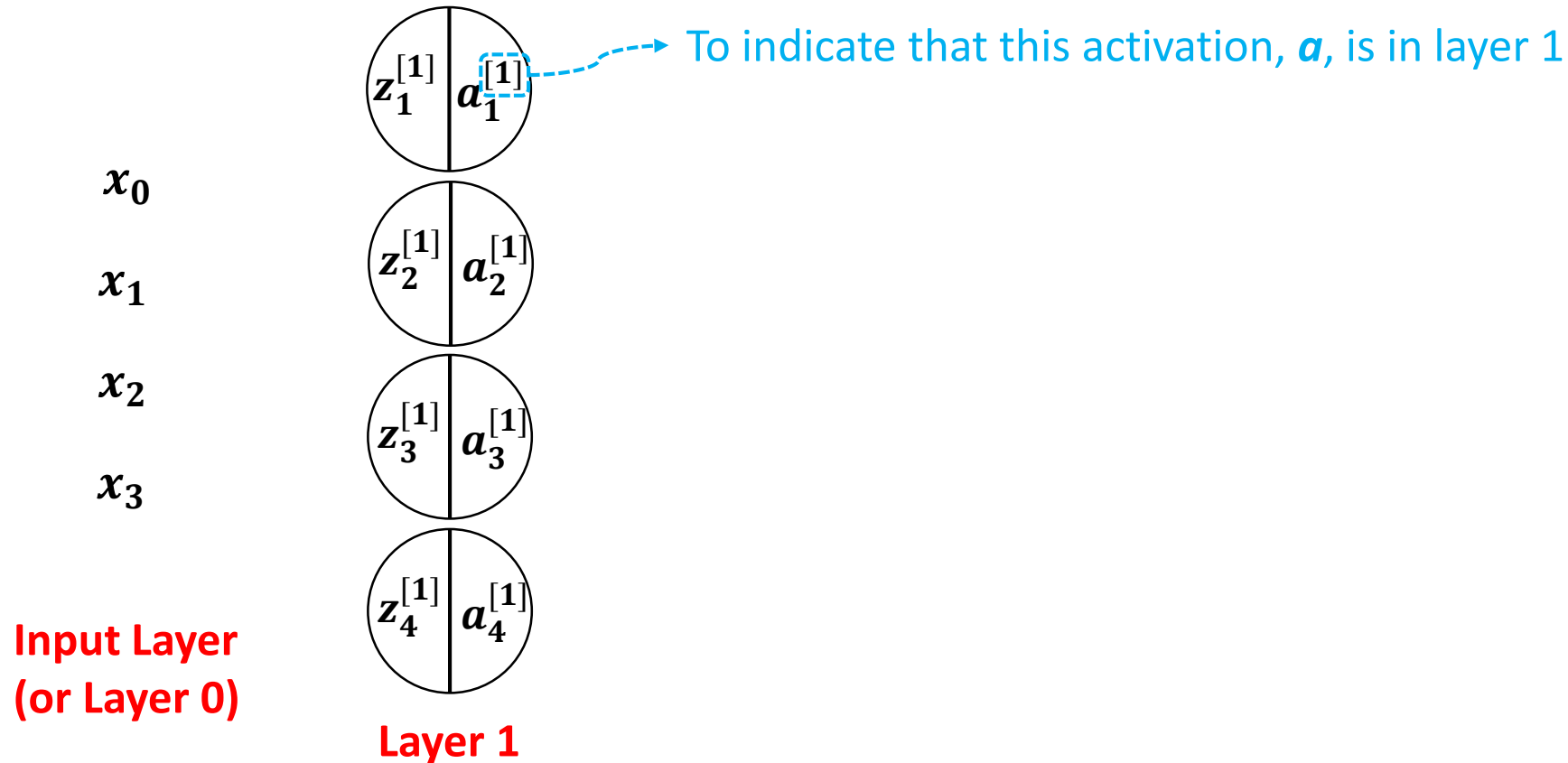Example from to Eric Postma via Jason Eisner

# playground.tensorflow.org

# Neural Networks

- We can construct a network of neurons (i.e., a neural network) with as many *layers,* and neurons in any layer, as needed

$x_0$

$x_1$

$x_2$

$x_3$

**Input Layer (or Layer 0)**

$z_1^{[1]}$ $a_1^{[1]}$

To indicate that this activation, *a*, is in layer 1

$z_2^{[1]}$ $a_2^{[1]}$
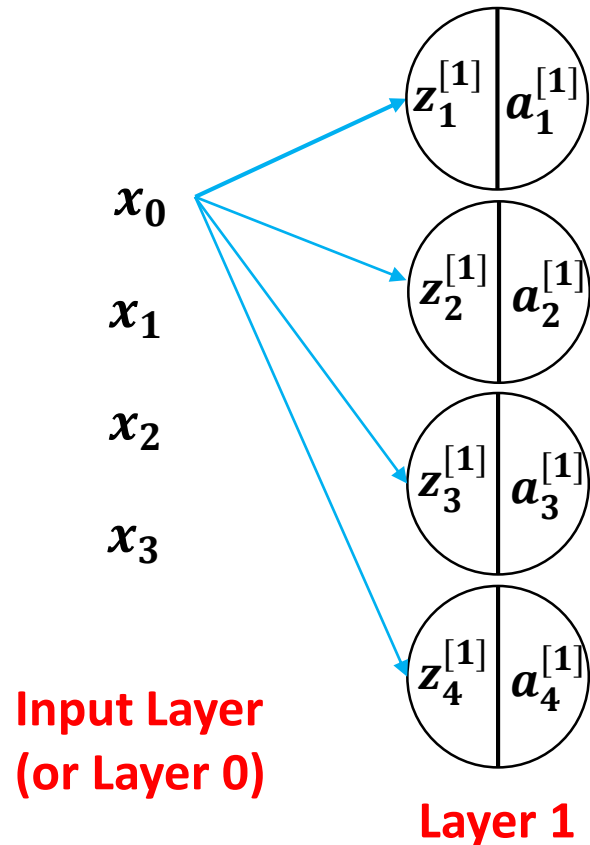
$z_3^{[1]}$ $a_3^{[1]}$

$z_4^{[1]}$ $a_4^{[1]}$

**Layer 1**

# Neural Networks

- We can construct a network of neurons (i.e., a neural network) with as many *layers,* and neurons in any layer, as needed
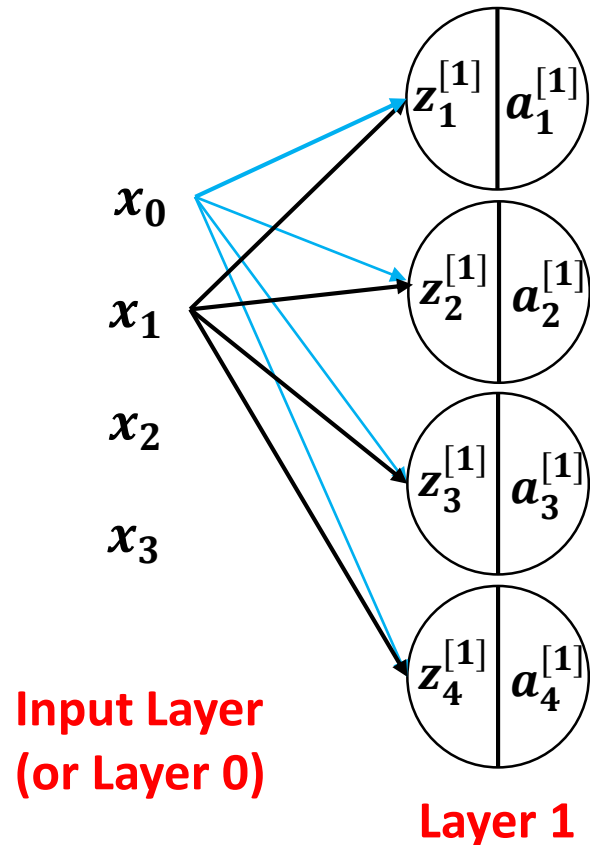
$x_0$

$x_1$

$x_2$

$x_3$

$z_1^{[1]}$ | $a_1^{[1]}$

$z_2^{[1]}$ | $a_2^{[1]}$

$z_3^{[1]}$ | $a_3^{[1]}$

$z_4^{[1]}$ | $a_4^{[1]}$

**Input Layer
(or Layer 0)**

**Layer 1**

# Neural Networks

- We can construct a network of neurons (i.e., a neural network) with as many *layers,* and neurons in any layer, as needed



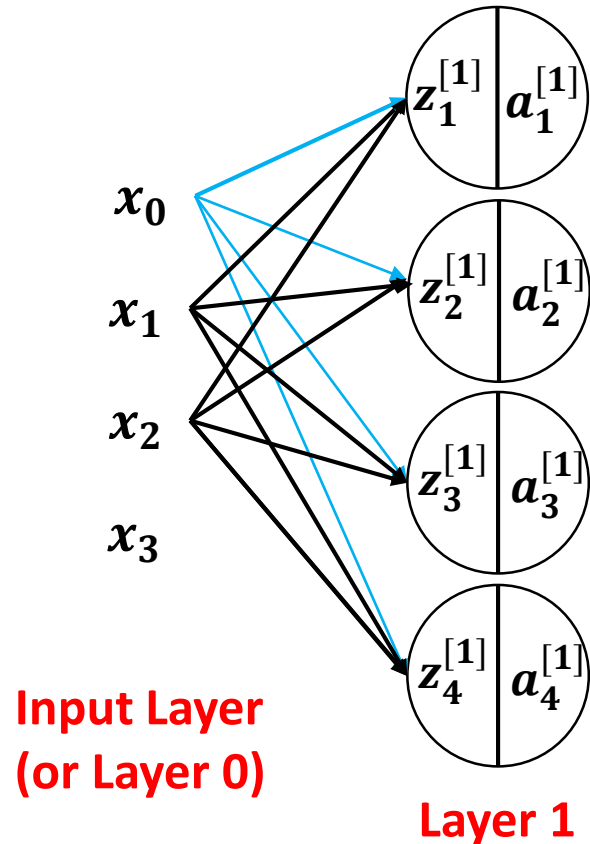$x_0$

$x_1$

$x_2$

$x_3$

$z_1^{[1]}$ | $a_1^{[1]}$

$z_2^{[1]}$ | $a_2^{[1]}$

$z_3^{[1]}$ | $a_3^{[1]}$

$z_4^{[1]}$ | $a_4^{[1]}$

**Input Layer
(or Layer 0)**

**Layer 1**

# Neural Networks

- We can construct a network of neurons (i.e., a neural network) with as many *layers,* and neurons in any layer, as needed
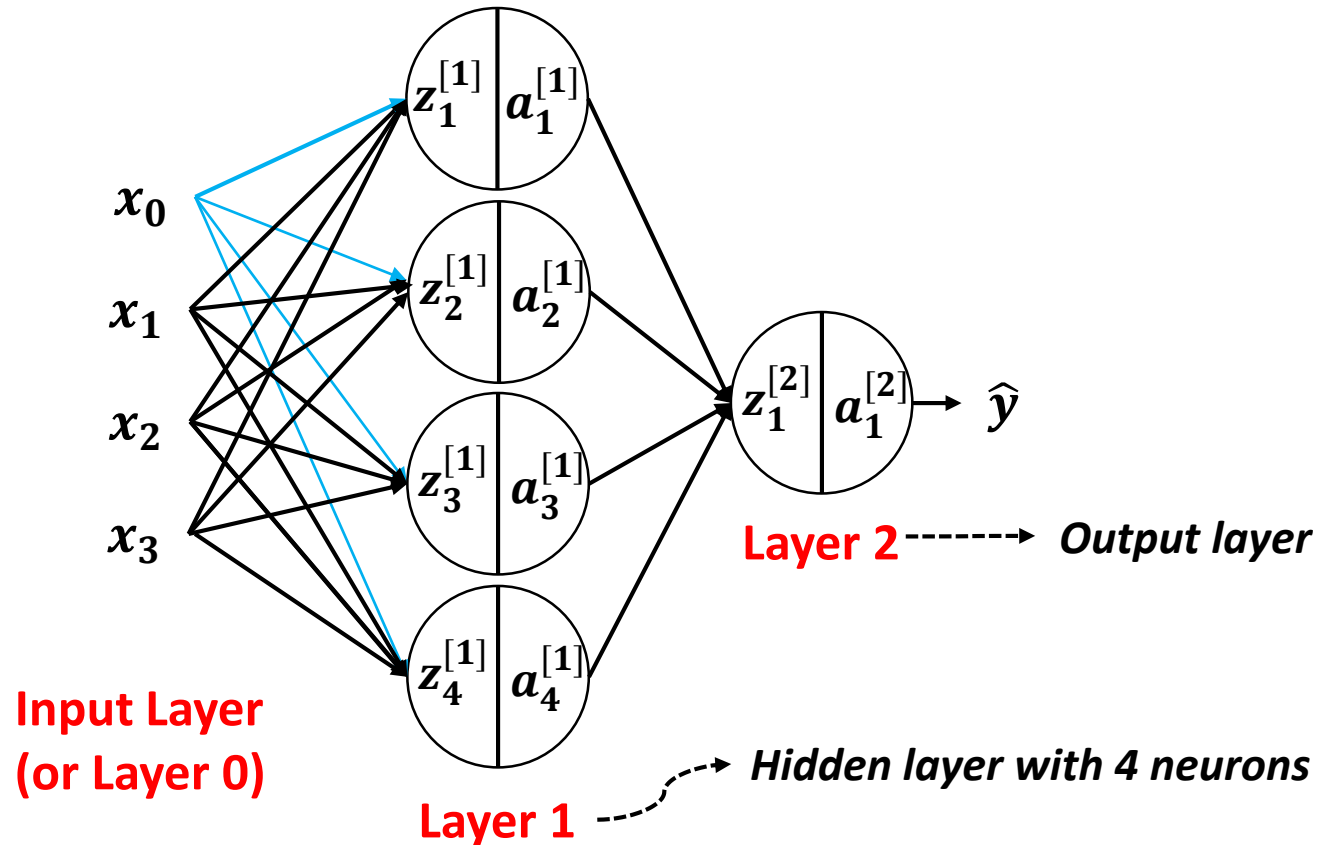


$x_0$

$x_1$

$x_2$

$x_3$

$z_1^{[1]}$ | $a_1^{[1]}$

$z_2^{[1]}$ | $a_2^{[1]}$

$z_3^{[1]}$ | $a_3^{[1]}$

$z_4^{[1]}$ | $a_4^{[1]}$

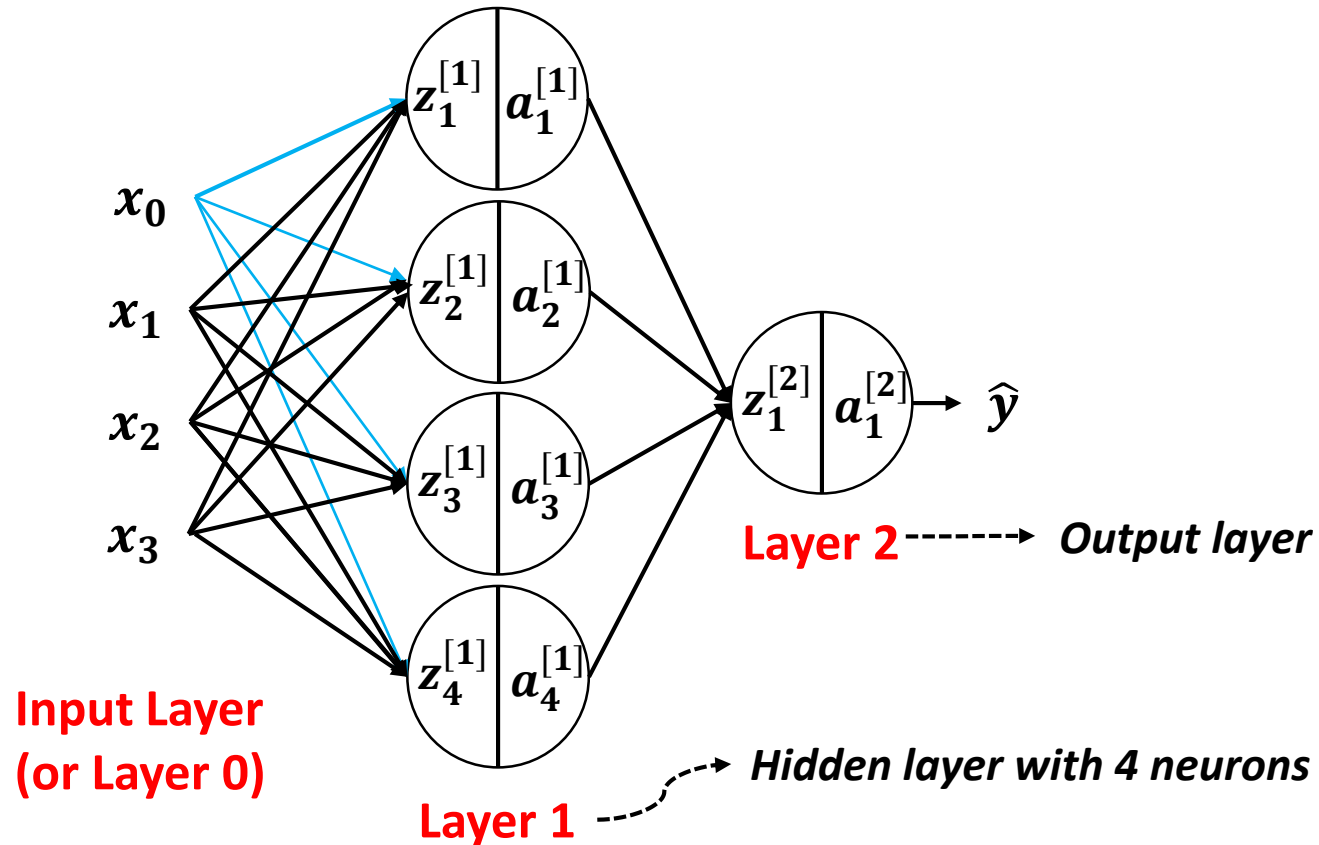**Input Layer (or Layer 0)**

**Layer 1**

# Neural Networks

- We can construct a network of neurons (i.e., a neural network) with as many *layers,* and neurons in any layer, as needed



By convention, this neural network is said to have 2 layers (and not 3) since the input layer is typically not counted!
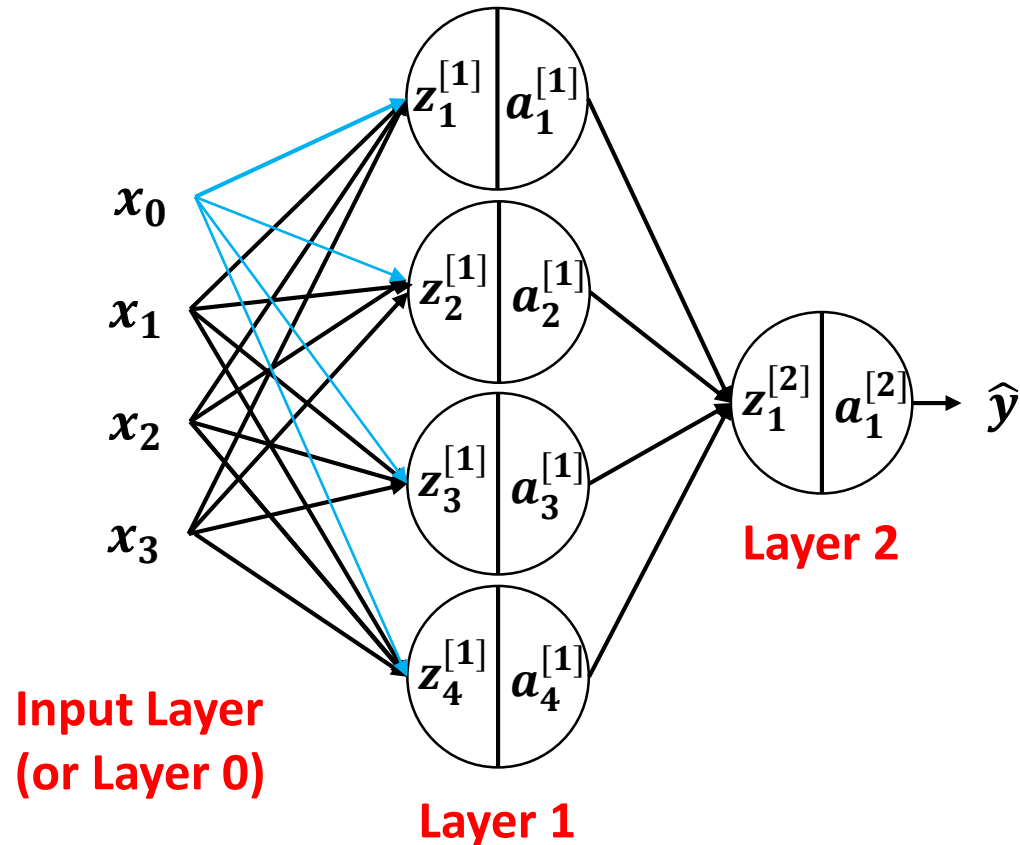
# Neural Networks

- We can construct a network of neurons (i.e., a neural network) with as many *layers,* and neurons in any layer, as needed



$x_0$

$x_1$

$x_2$

$x_3$

$z_1^{[1]}$ | $a_1^{[1]}$

$z_2^{[1]}$ | $a_2^{[1]}$

$z_3^{[1]}$ | $a_3^{[1]}$

$z_4^{[1]}$ | $a_4^{[1]}$

$z_1^{[2]}$ | $a_1^{[2]}$ → $\hat{y}$

**Layer 2** - - - - - → *Output layer*

**Input Layer (or Layer 0)**

- - → *Hidden layer with 4 neurons*

**Layer 1** - - -

Also, the more layers we add, the ***deeper*** the neural network becomes, giving rise to the concept of ***deep learning***!

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us *vectorize* (represent in vectors & matrices) the input and the variables involved
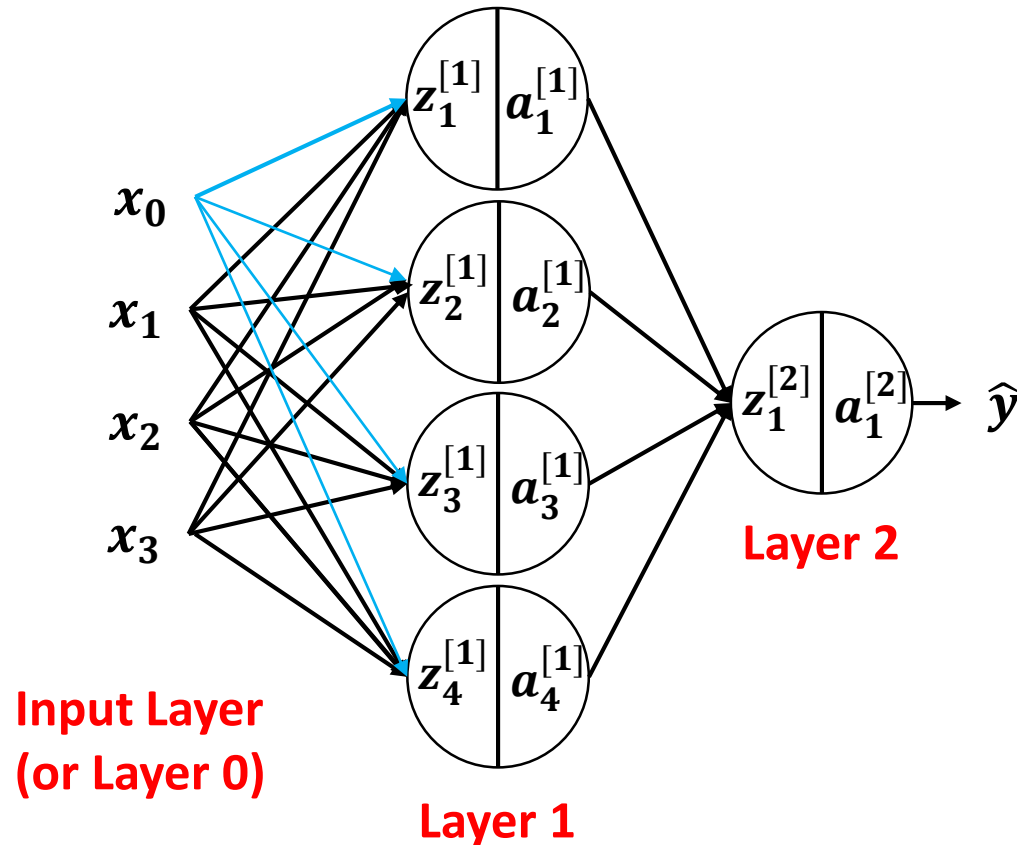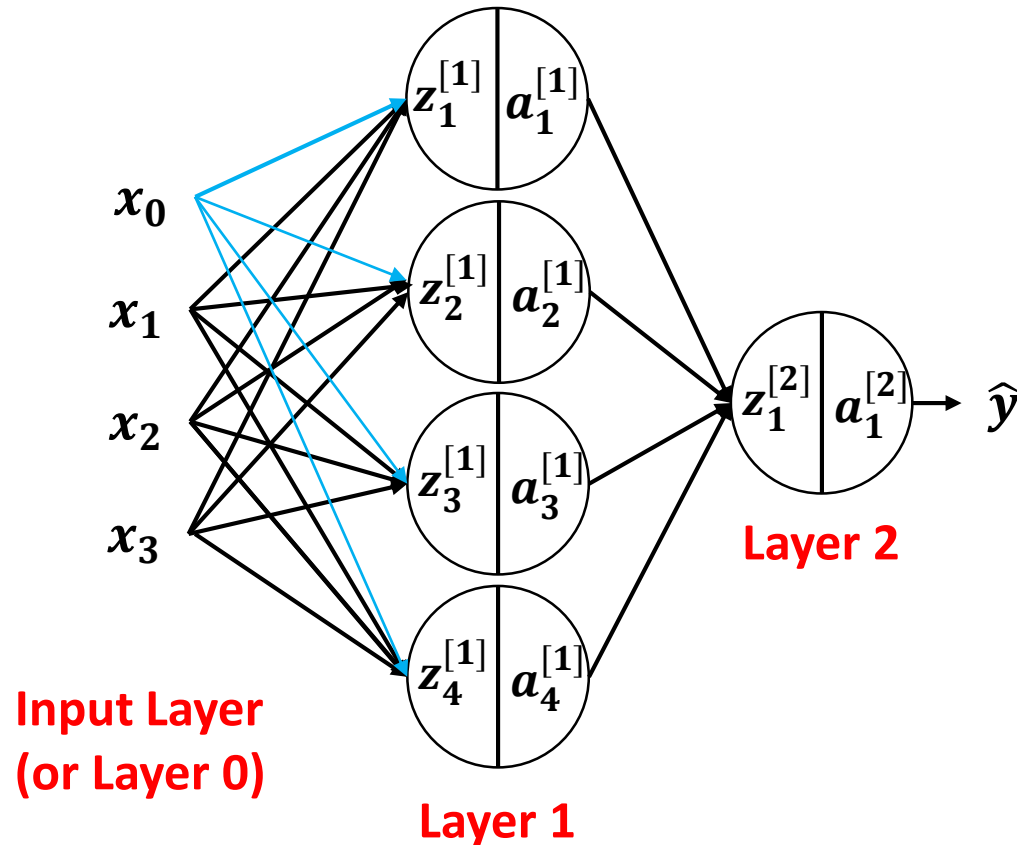


$$z^{[1]} = w^{[1]^T}x + b^{[1]}$$

$$= \begin{bmatrix} w_{11}^{[1]} & w_{21}^{[1]} & w_{31}^{[1]} \\ w_{12}^{[1]} & w_{22}^{[1]} & w_{32}^{[1]} \\ w_{13}^{[1]} & w_{23}^{[1]} & w_{33}^{[1]} \\ w_{14}^{[1]} & w_{24}^{[1]} & w_{34}^{[1]} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}$$

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us **vectorize** (represent in vectors & matrices) the input and the variables involved
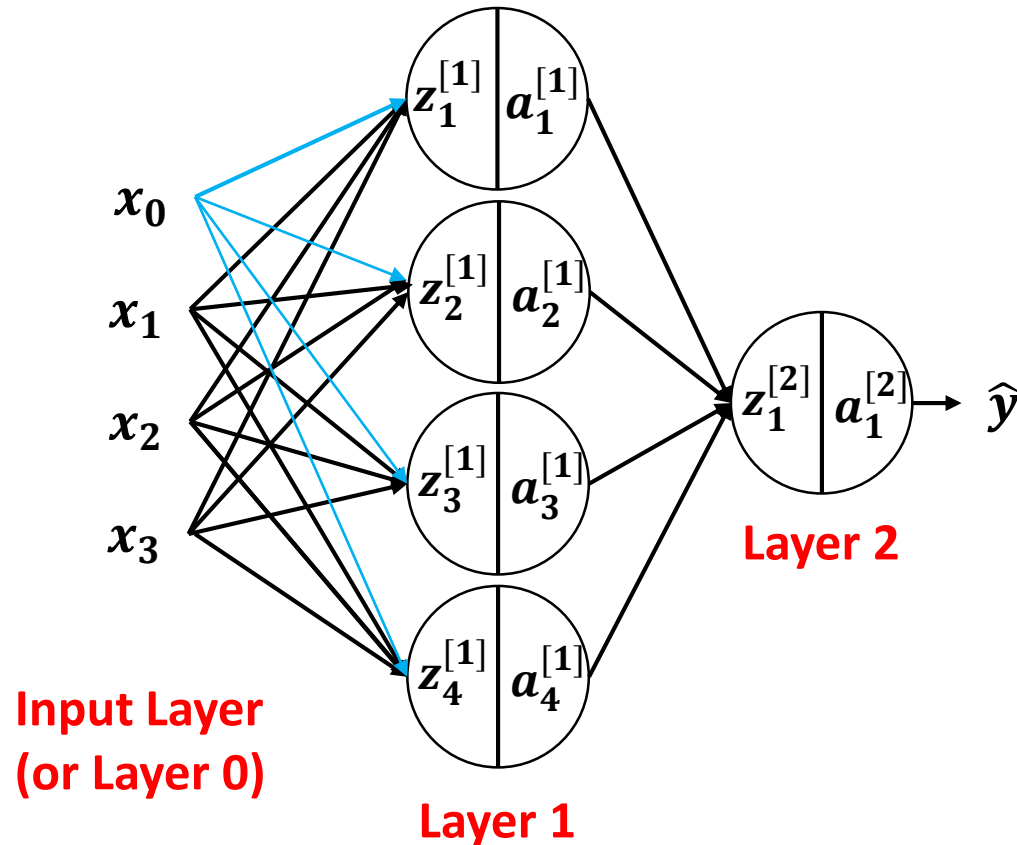


$$z^{[1]} = w^{[1]^T} x + b^{[1]}$$

$$= \begin{bmatrix} w_{11}^{[1]} x_1 + w_{21}^{[1]} x_2 + w_{31}^{[1]} x_3 + b_1^{[1]} \\ w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2 + w_{32}^{[1]} x_3 + b_2^{[1]} \\ w_{13}^{[1]} x_1 + w_{23}^{[1]} x_2 + w_{33}^{[1]} x_3 + b_3^{[1]} \\ w_{14}^{[1]} x_1 + w_{24}^{[1]} x_2 + w_{34}^{[1]} x_3 + b_4^{[1]} \end{bmatrix}$$

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

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us *vectorize* (represent in vectors & matrices) the input and the variables involved
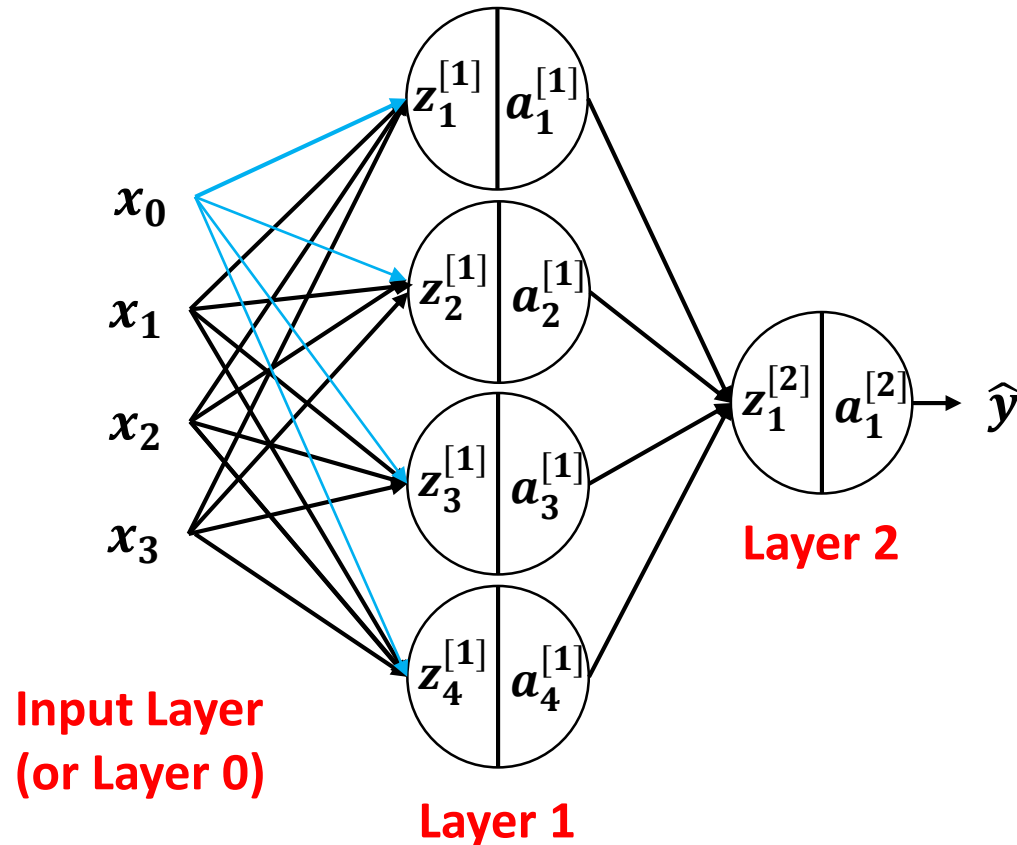


$$z^{[1]} = w^{[1]^T}x + b^{[1]}$$

$$= \begin{bmatrix} w_{11}^{[1]}x_1 + w_{21}^{[1]}x_2 + w_{31}^{[1]}x_3 + b_1^{[1]} \\ w_{12}^{[1]}x_1 + w_{22}^{[1]}x_2 + w_{32}^{[1]}x_3 + b_2^{[1]} \\ w_{13}^{[1]}x_1 + w_{23}^{[1]}x_2 + w_{33}^{[1]}x_3 + b_3^{[1]} \\ w_{14}^{[1]}x_1 + w_{24}^{[1]}x_2 + w_{34}^{[1]}x_3 + b_4^{[1]} \end{bmatrix}$$

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

$$z_1^{[1]}$$

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us *vectorize* (represent in vectors & matrices) the input and the variables involved
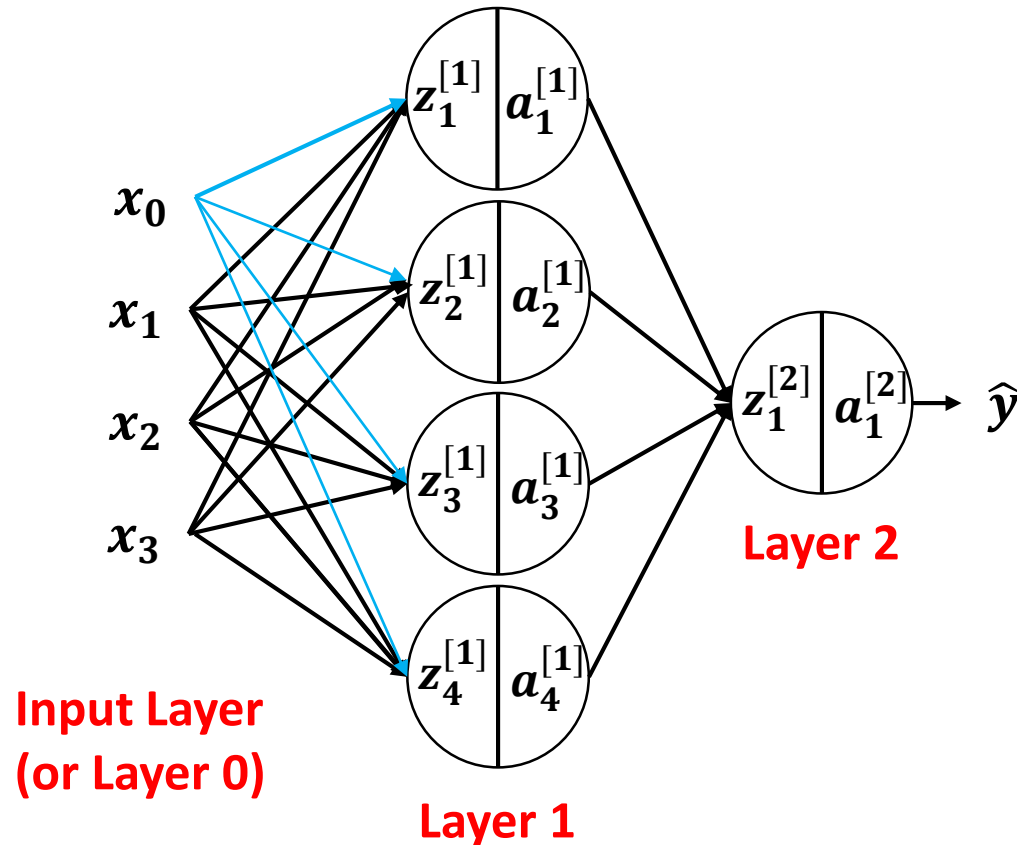


$$z^{[1]} = w^{[1]^T} x + b^{[1]}$$

$$= \begin{bmatrix} w_{11}^{[1]} x_1 + w_{21}^{[1]} x_2 + w_{31}^{[1]} x_3 + b_1^{[1]} \\ w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2 + w_{32}^{[1]} x_3 + b_2^{[1]} \\ w_{13}^{[1]} x_1 + w_{23}^{[1]} x_2 + w_{33}^{[1]} x_3 + b_3^{[1]} \\ w_{14}^{[1]} x_1 + w_{24}^{[1]} x_2 + w_{34}^{[1]} x_3 + b_4^{[1]} \end{bmatrix}$$

$$a^{[1]} = \sigma(z^{[1]}) \qquad\qquad a_1^{[1]} = \sigma(z_1^{[1]})$$

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us **vectorize** (represent in vectors & matrices) the input and the variables involved
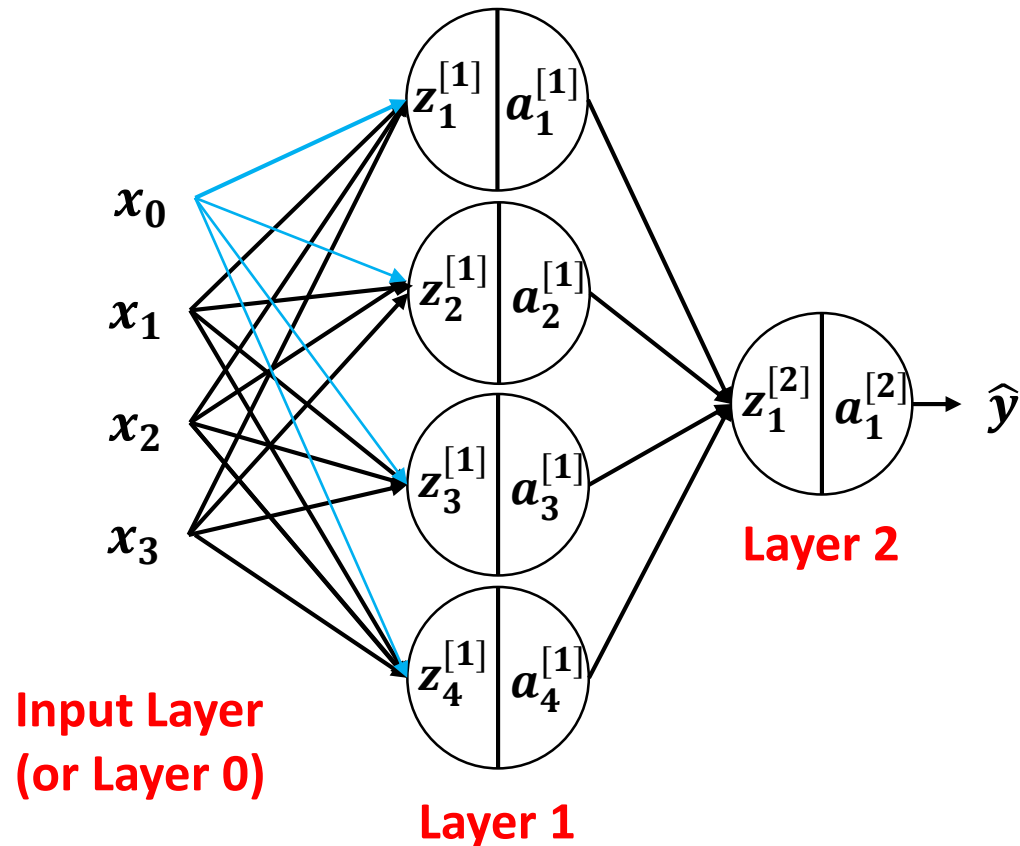


$$z^{[1]} = w^{[1]^T} x + b^{[1]}$$

$$= \begin{bmatrix} w_{11}^{[1]} x_1 + w_{21}^{[1]} x_2 + w_{31}^{[1]} x_3 + b_1^{[1]} \\ w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2 + w_{32}^{[1]} x_3 + b_2^{[1]} \\ w_{13}^{[1]} x_1 + w_{23}^{[1]} x_2 + w_{33}^{[1]} x_3 + b_3^{[1]} \\ w_{14}^{[1]} x_1 + w_{24}^{[1]} x_2 + w_{34}^{[1]} x_3 + b_4^{[1]} \end{bmatrix}$$

$$z_2^{[1]}$$

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

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us **vectorize** (represent in vectors & matrices) the input and the variables involved



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

$$= \begin{bmatrix} w_{11}^{[1]} x_1 + w_{21}^{[1]} x_2 + w_{31}^{[1]} x_3 + b_1^{[1]} \\ w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2 + w_{32}^{[1]} x_3 + b_2^{[1]} \\ w_{13}^{[1]} x_1 + w_{23}^{[1]} x_2 + w_{33}^{[1]} x_3 + b_3^{[1]} \\ w_{14}^{[1]} x_1 + w_{24}^{[1]} x_2 + w_{34}^{[1]} x_3 + b_4^{[1]} \end{bmatrix}$$
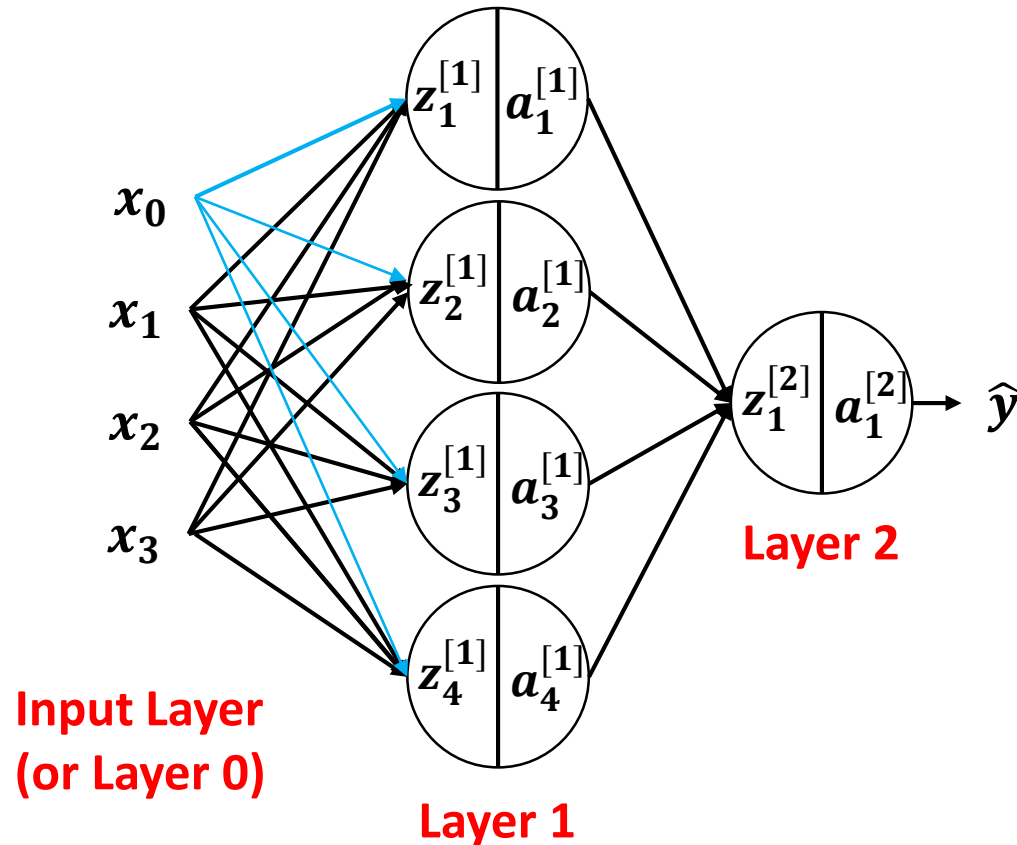
$$a^{[1]} = \sigma(z^{[1]})$$

$$a_2^{[1]} = \sigma(z_2^{[1]})$$

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us *vectorize* (represent in vectors & matrices) the input and the variables involved



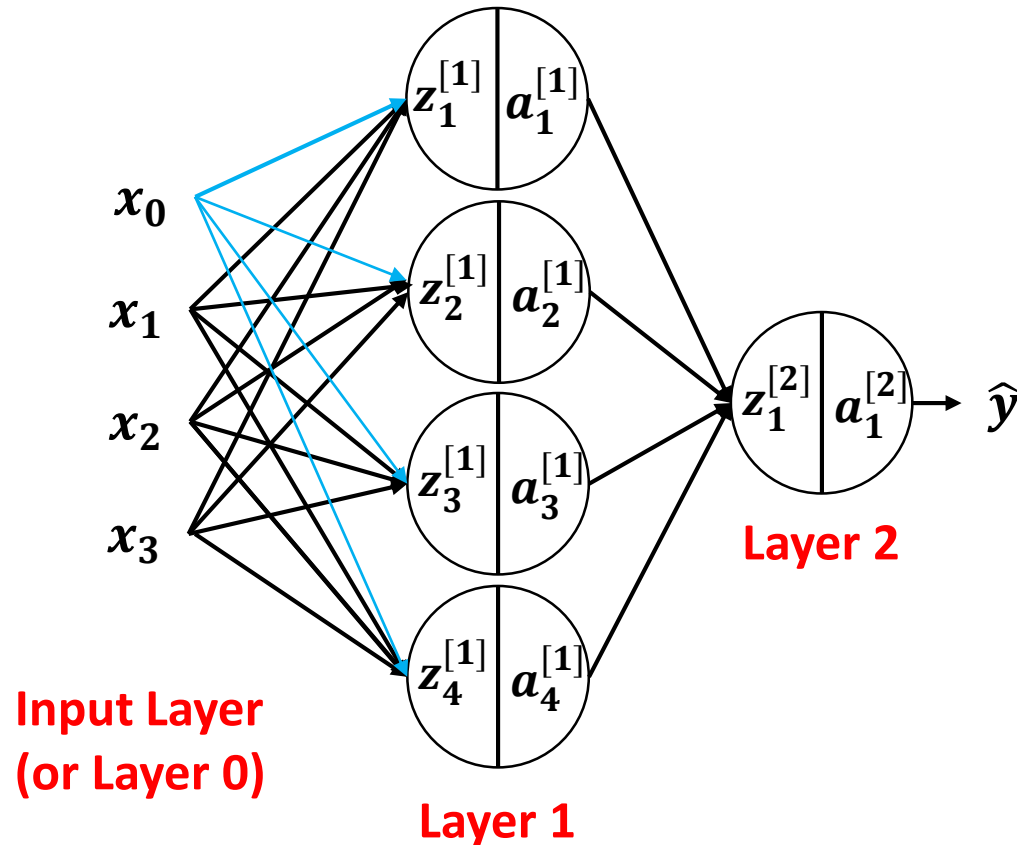$$\boldsymbol{z}^{[1]} = \boldsymbol{w}^{[1]^T}\boldsymbol{x} + \boldsymbol{b}^{[1]}$$

$$= \begin{bmatrix} w_{11}^{[1]}x_1 + w_{21}^{[1]}x_2 + w_{31}^{[1]}x_3 + b_1^{[1]} \\ w_{12}^{[1]}x_1 + w_{22}^{[1]}x_2 + w_{32}^{[1]}x_3 + b_2^{[1]} \\ w_{13}^{[1]}x_1 + w_{23}^{[1]}x_2 + w_{33}^{[1]}x_3 + b_3^{[1]} \\ w_{14}^{[1]}x_1 + w_{24}^{[1]}x_2 + w_{34}^{[1]}x_3 + b_4^{[1]} \end{bmatrix}$$

$$\boldsymbol{a}^{[1]} = \boldsymbol{\sigma}(\boldsymbol{z}^{[1]})$$

$$z_3^{[1]}$$

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us **_vectorize_** (represent in vectors & matrices) the input and the variables involved
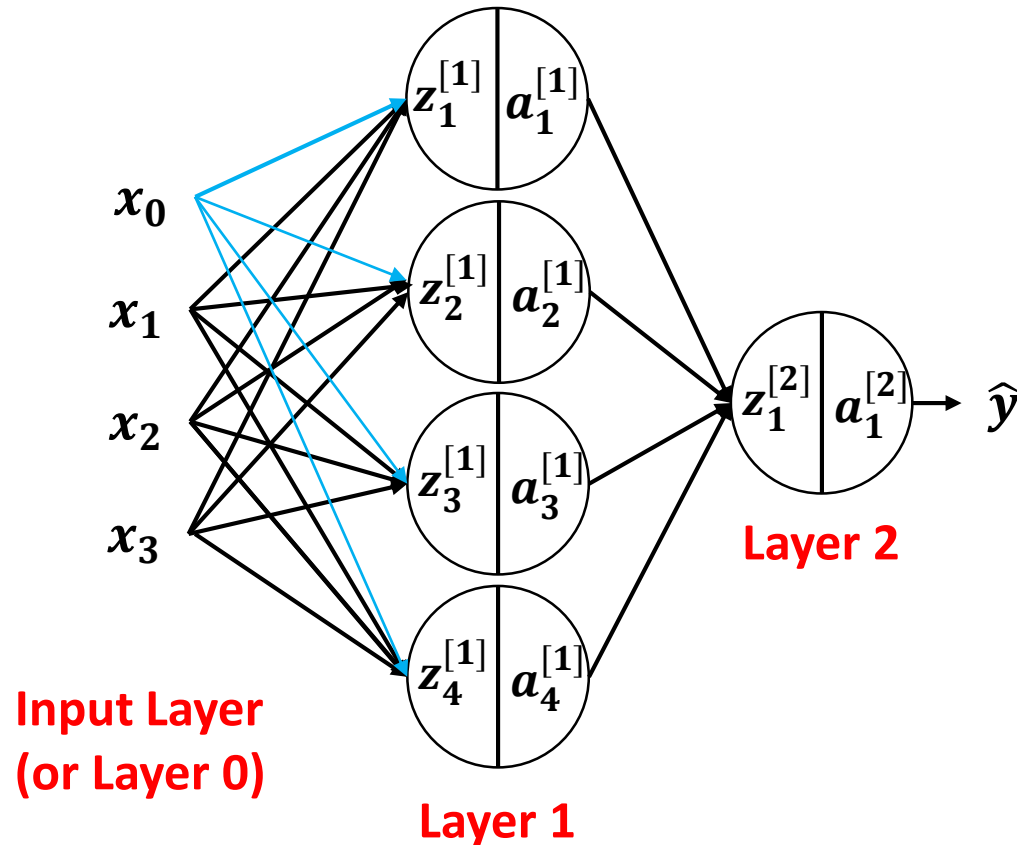


$$z^{[1]} = w^{[1]^T}x + b^{[1]}$$

$$= \begin{bmatrix} w_{11}^{[1]}x_1 + w_{21}^{[1]}x_2 + w_{31}^{[1]}x_3 + b_1^{[1]} \\ w_{12}^{[1]}x_1 + w_{22}^{[1]}x_2 + w_{32}^{[1]}x_3 + b_2^{[1]} \\ w_{13}^{[1]}x_1 + w_{23}^{[1]}x_2 + w_{33}^{[1]}x_3 + b_3^{[1]} \\ w_{14}^{[1]}x_1 + w_{24}^{[1]}x_2 + w_{34}^{[1]}x_3 + b_4^{[1]} \end{bmatrix}$$

$$a^{[1]} = \sigma(z^{[1]}) \qquad a_3^{[1]} = \sigma(z_3^{[1]})$$

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us *vectorize* (represent in vectors & matrices) the input and the variables involved
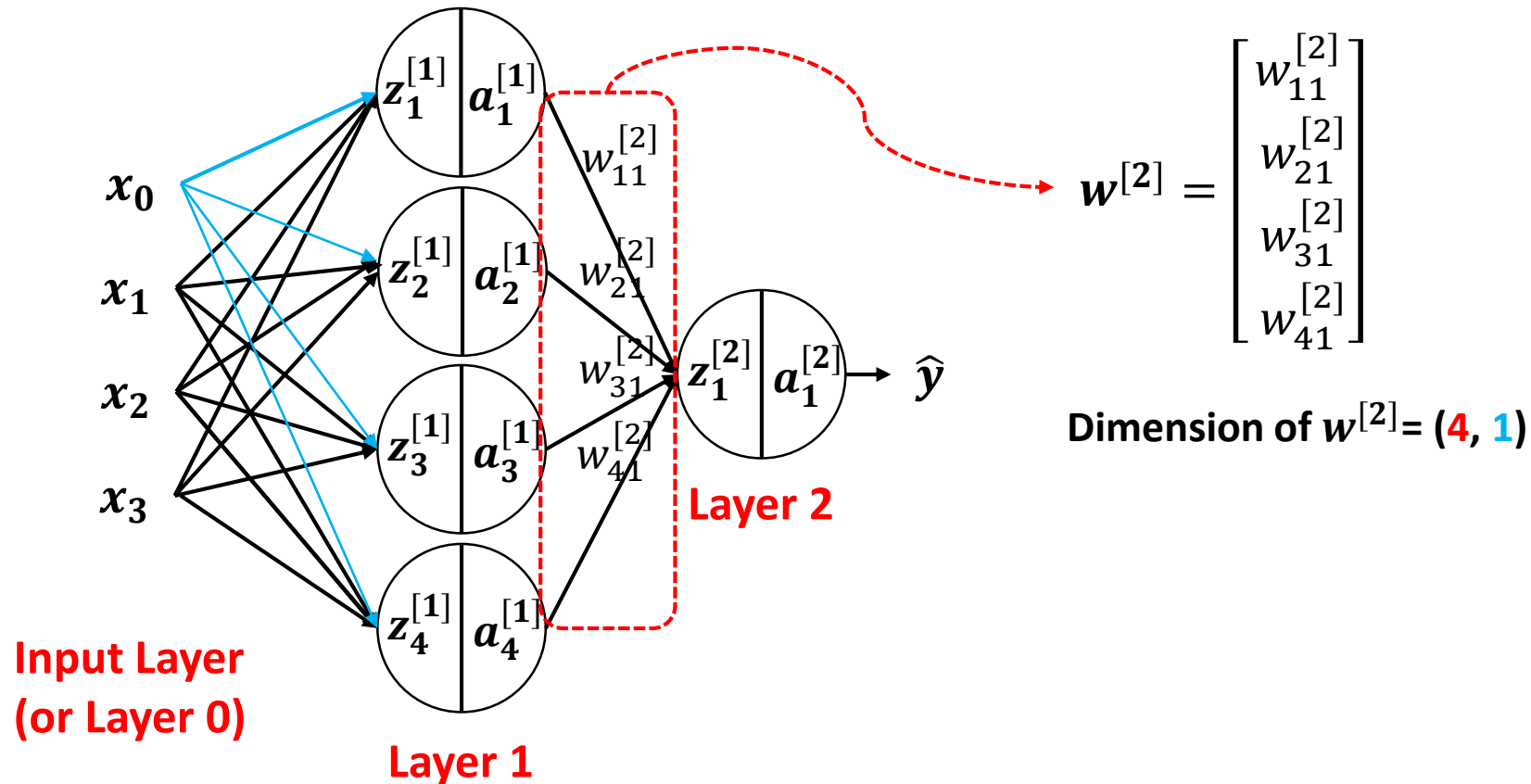


$$z^{[1]} = w^{[1]^T} x + b^{[1]}$$

$$= \begin{bmatrix} w_{11}^{[1]} x_1 + w_{21}^{[1]} x_2 + w_{31}^{[1]} x_3 + b_1^{[1]} \\ w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2 + w_{32}^{[1]} x_3 + b_2^{[1]} \\ w_{13}^{[1]} x_1 + w_{23}^{[1]} x_2 + w_{33}^{[1]} x_3 + b_3^{[1]} \\ w_{14}^{[1]} x_1 + w_{24}^{[1]} x_2 + w_{34}^{[1]} x_3 + b_4^{[1]} \end{bmatrix}$$

$$z_4^{[1]}$$
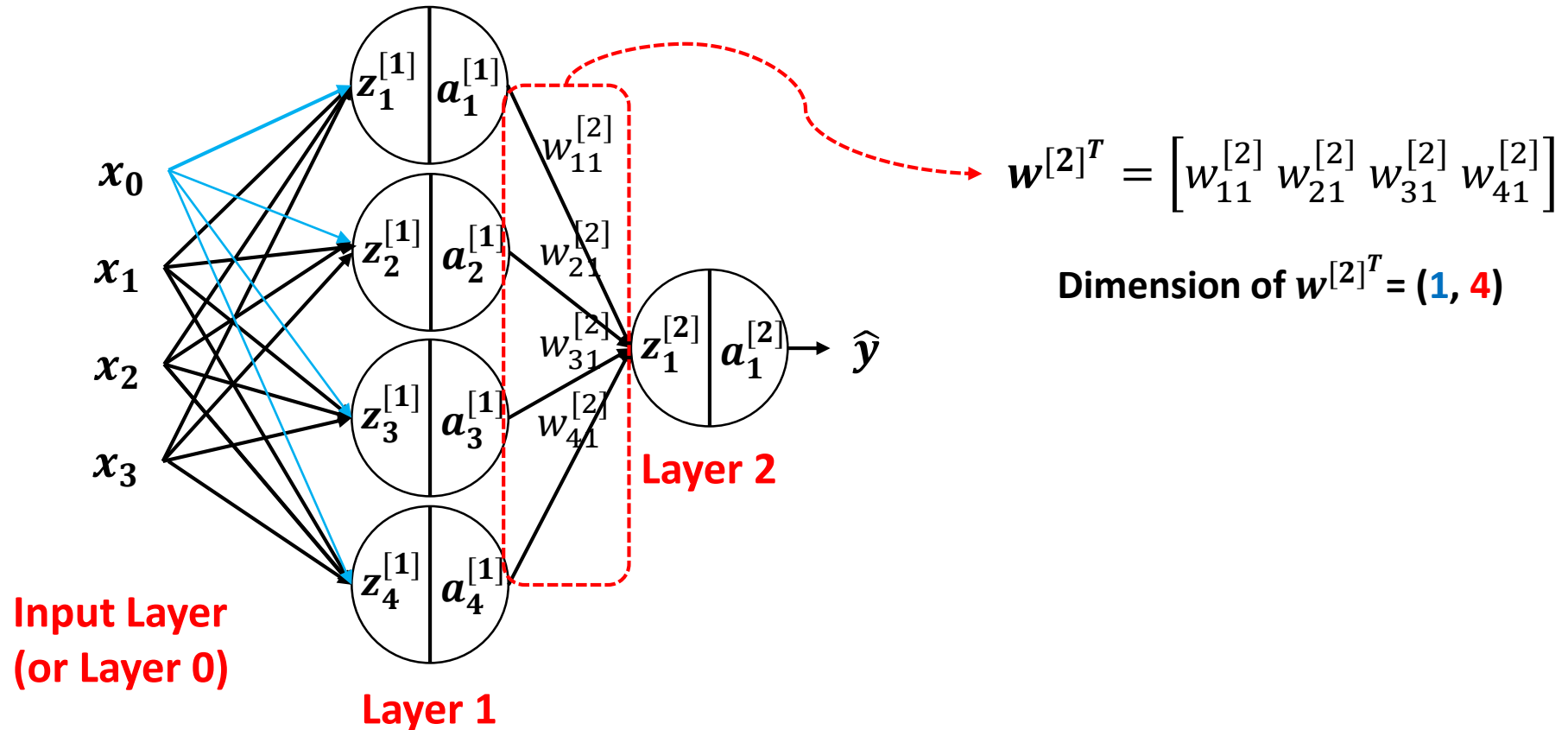
$$a^{[1]} = \sigma(z^{[1]})$$

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us **_vectorize_** (represent in vectors & matrices) the input and the variables involved



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

$$= \begin{bmatrix} w_{11}^{[1]} x_1 + w_{21}^{[1]} x_2 + w_{31}^{[1]} x_3 + b_1^{[1]} \\ w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2 + w_{32}^{[1]} x_3 + b_2^{[1]} \\ w_{13}^{[1]} x_1 + w_{23}^{[1]} x_2 + w_{33}^{[1]} x_3 + b_3^{[1]} \\ w_{14}^{[1]} x_1 + w_{24}^{[1]} x_2 + w_{34}^{[1]} x_3 + b_4^{[1]} \end{bmatrix}$$

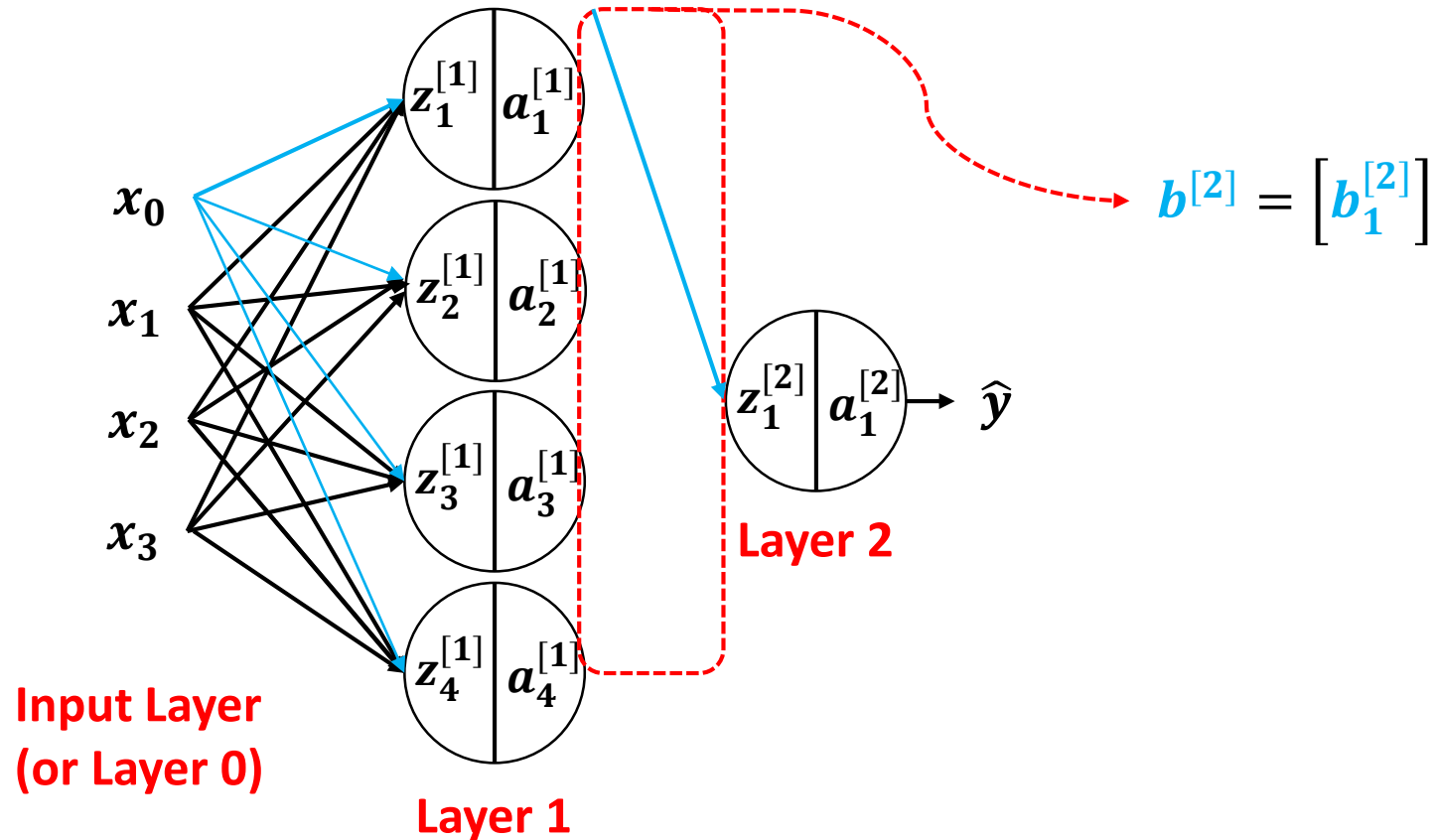$$a^{[1]} = \sigma(z^{[1]}) \qquad\qquad a_4^{[1]} = \sigma(z_4^{[1]})$$

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us *vectorize* (represent in vectors & matrices) the input and the variables involved
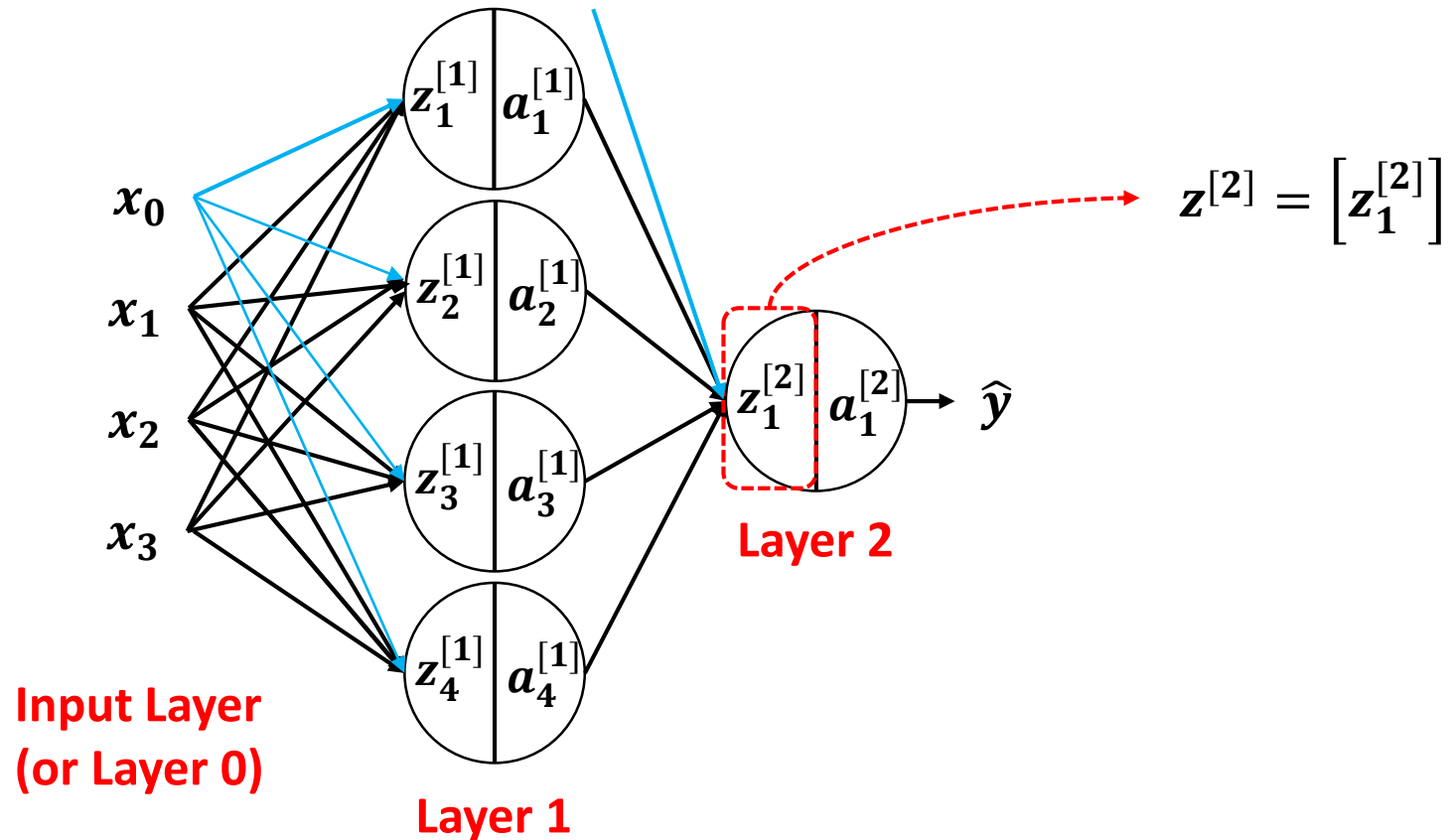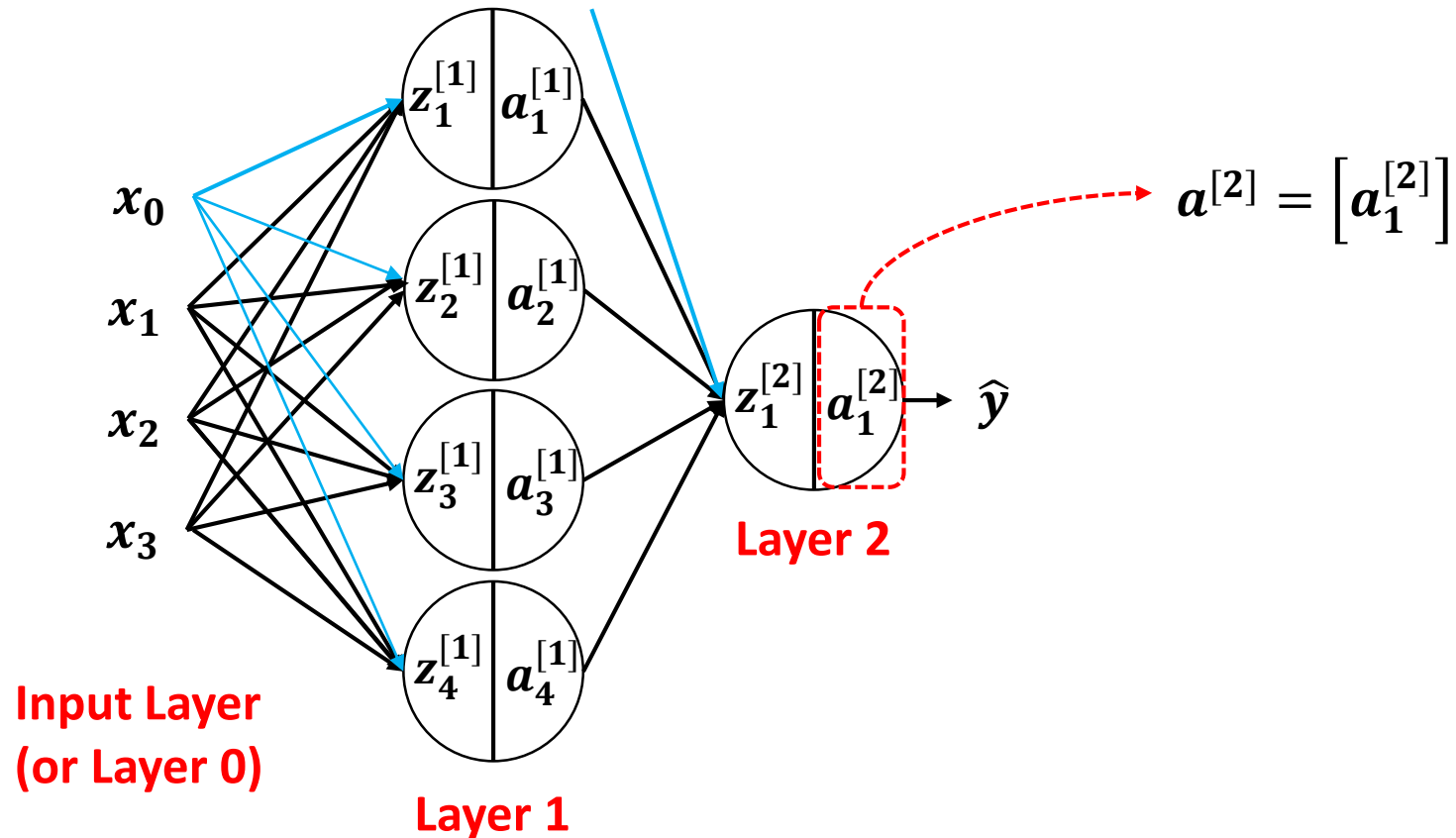
# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us **_vectorize_** (represent in vectors & matrices) the input and the variables involved



$$w^{[2]^T} = \begin{bmatrix} w_{11}^{[2]} & w_{21}^{[2]} & w_{31}^{[2]} & w_{41}^{[2]} \end{bmatrix}$$

**Dimension of $w^{[2]^T}$ = (1, 4)**

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us *vectorize* (represent in vectors & matrices) the input and the variables involved
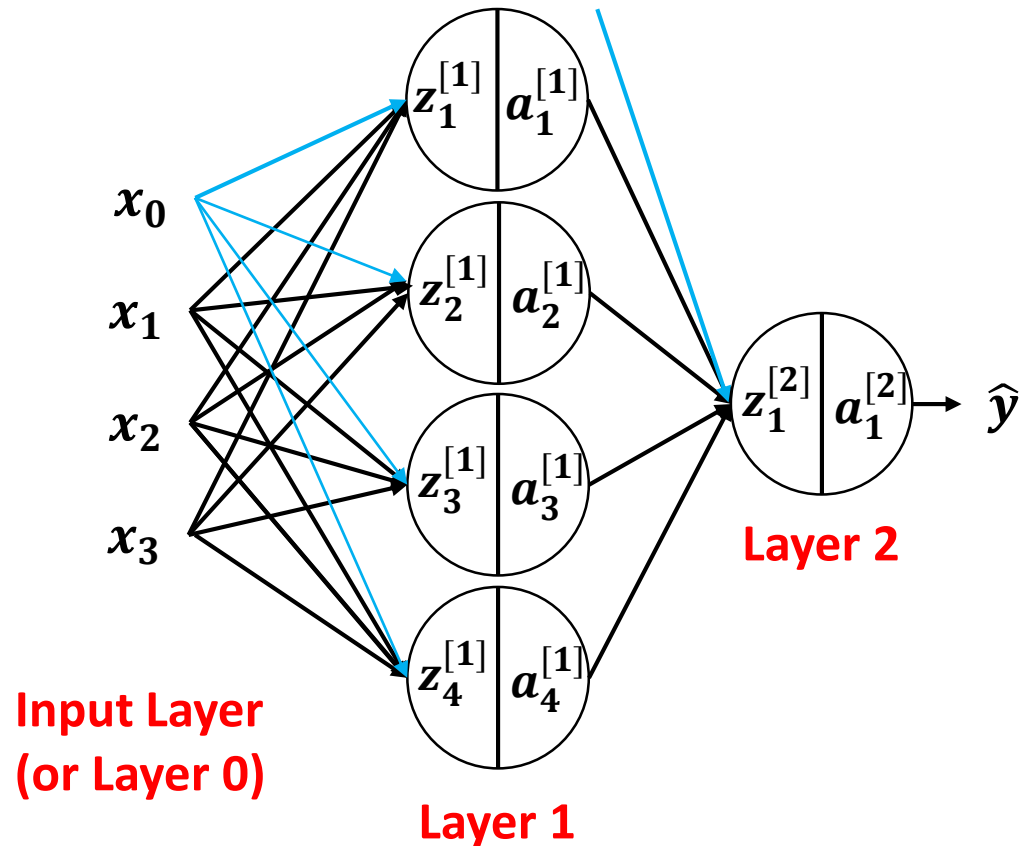
# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us **vectorize** (represent in vectors & matrices) the input and the variables involved



$$z^{[2]} = \left[ z_1^{[2]} \right]$$

**Layer 2**

**Input Layer (or Layer 0)**

**Layer 1**

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us *vectorize* (represent in vectors & matrices) the input and the variables involved
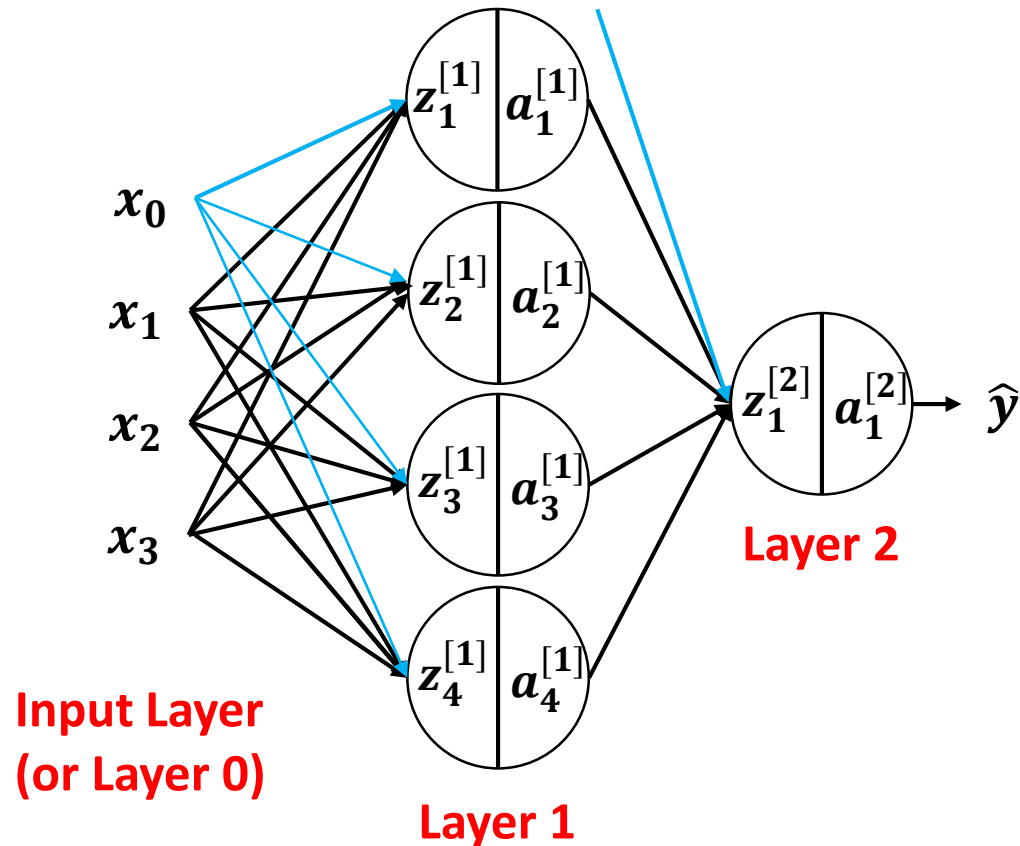
# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us **_vectorize_** (represent in vectors & matrices) the input and the variables involved



$$z^{[2]} = w^{[2]^T}x + b^{[2]}$$

$$= w^{[2]^T}a^{[1]} + b^{[2]}$$

$$= \begin{bmatrix} w^{[2]}_{11} & w^{[2]}_{21} & w^{[2]}_{31} & w^{[2]}_{41} \end{bmatrix} \begin{bmatrix} a^{[1]}_1 \\ a^{[1]}_2 \\ a^{[1]}_3 \\ a^{[1]}_4 \end{bmatrix} + \begin{bmatrix} b^{[2]}_1 \end{bmatrix}$$

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us *vectorize* (represent in vectors & matrices) the input and the variables involved
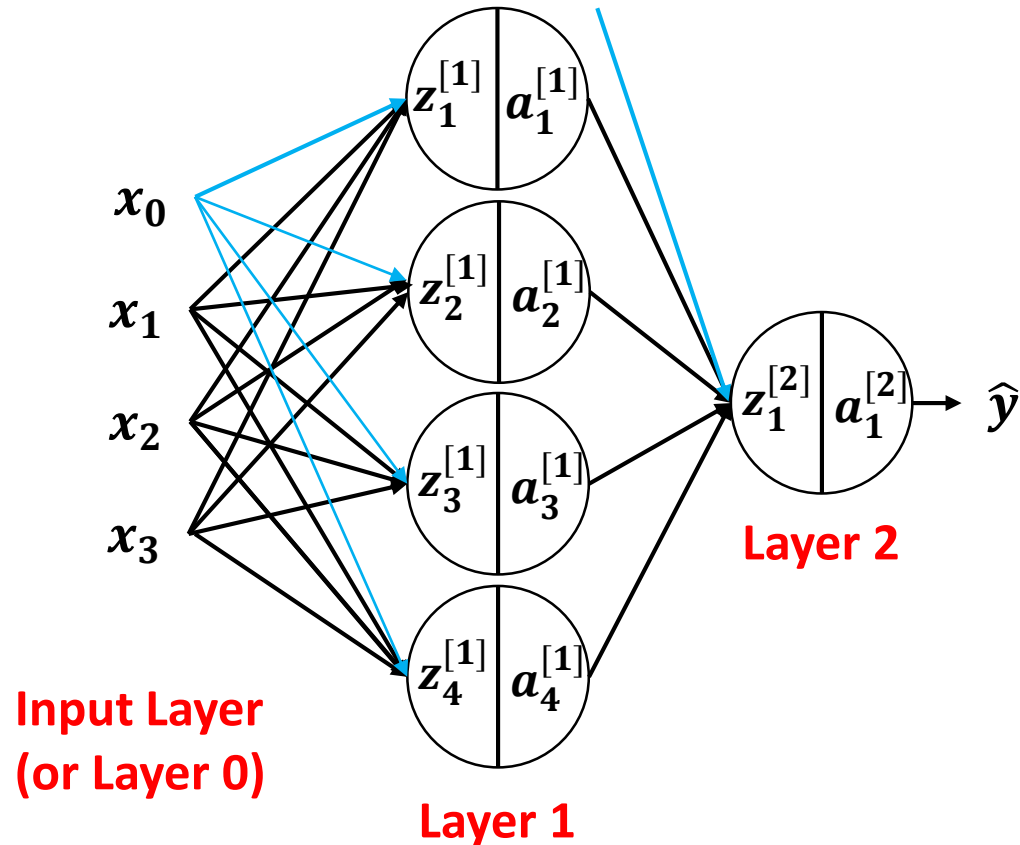


$$z^{[2]} = w^{[2]^T} x + b^{[2]}$$

$$= w^{[2]^T} a^{[1]} + b^{[2]}$$

$$= \left[ w_{11}^{[2]} a_1^{[1]} + w_{21}^{[2]} a_2^{[1]} + w_{31}^{[2]} a_3^{[1]} + w_{41}^{[2]} a_4^{[1]} \right] + \left[ b_1^{[2]} \right]$$

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us **_vectorize_** (represent in vectors & matrices) the input and the variables involved
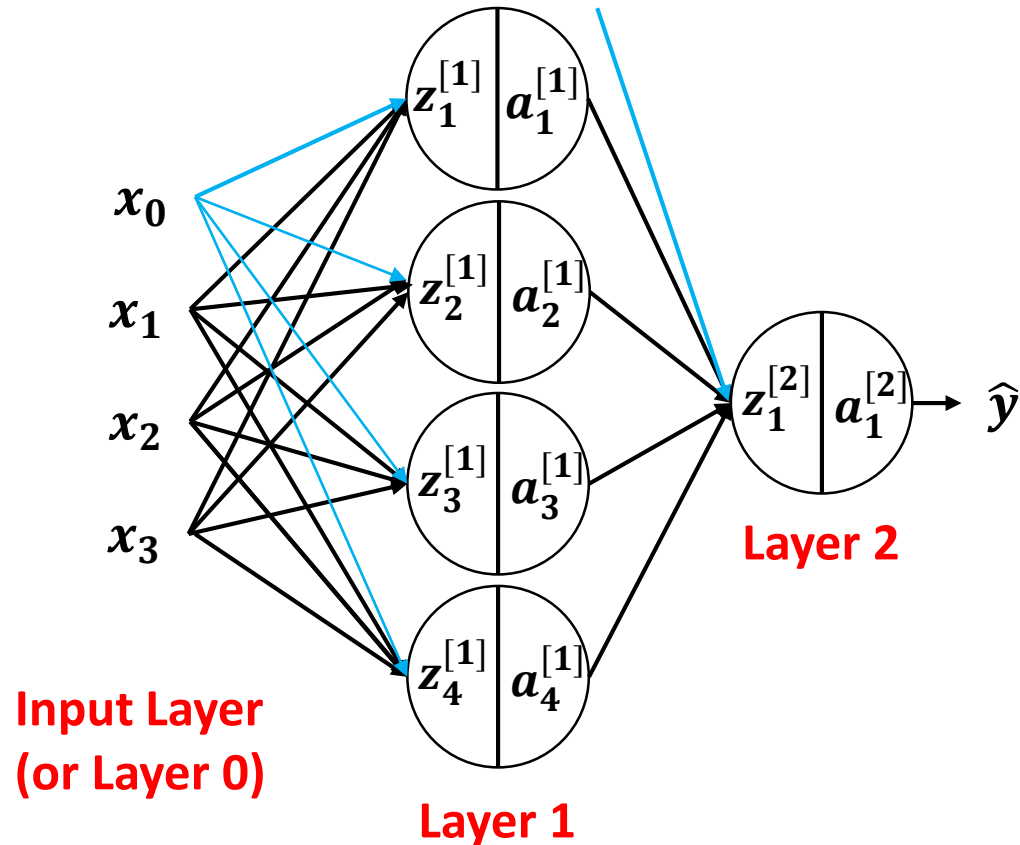


$$z^{[2]} = w^{[2]^T} x + b^{[2]}$$

$$= w^{[2]^T} a^{[1]} + b^{[2]}$$

$$= \left[ w_{11}^{[2]} a_1^{[1]} + w_{21}^{[2]} a_2^{[1]} + w_{31}^{[2]} a_3^{[1]} + w_{41}^{[2]} a_4^{[1]} + b_1^{[2]} \right]$$

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

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us *vectorize* (represent in vectors & matrices) the input and the variables involved



$$z^{[2]} = w^{[2]^T} x + b^{[2]}$$
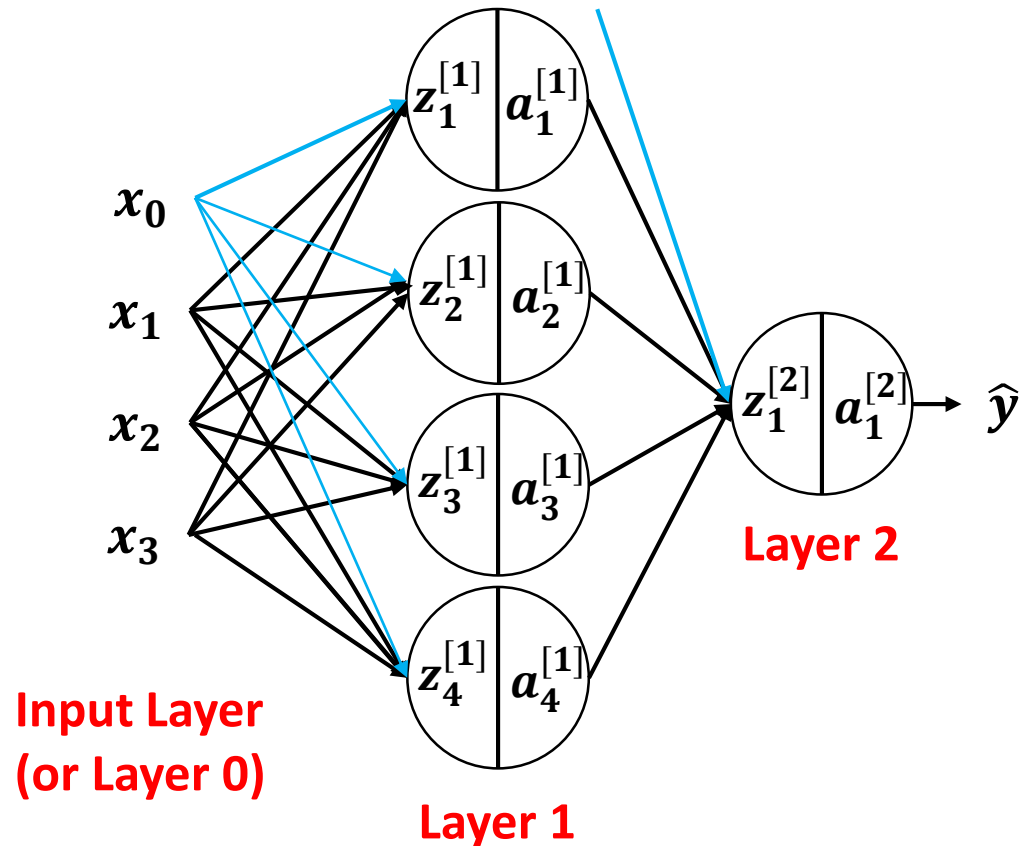
$$= w^{[2]^T} a^{[1]} + b^{[2]}$$

$$= \boxed{w_{11}^{[2]} a_1^{[1]} + w_{21}^{[2]} a_2^{[1]} + w_{31}^{[2]} a_3^{[1]} + w_{41}^{[2]} a_4^{[1]} + b_1^{[2]}}$$

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

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us *vectorize* (represent in vectors & matrices) the input and the variables involved



$$z^{[2]} = w^{[2]^T} x + b^{[2]}$$

$$= w^{[2]^T} a^{[1]} + b^{[2]}$$

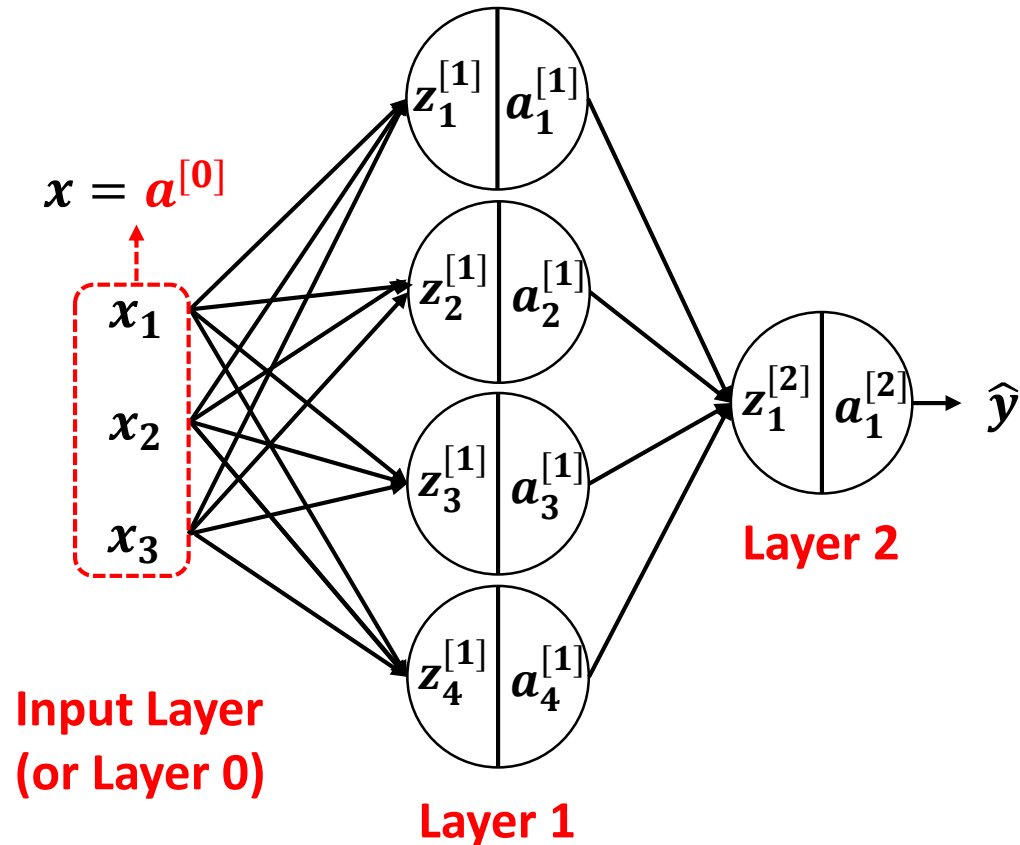$$= \boxed{w_{11}^{[2]} a_1^{[1]} + w_{21}^{[2]} a_2^{[1]} + w_{31}^{[2]} a_3^{[1]} + w_{41}^{[2]} a_4^{[1]} + b_1^{[2]}}$$

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

$$a_1^{[2]} = \sigma(z_1^{[2]})$$

**Input Layer (or Layer 0)**

**Layer 1**

**Layer 2**

# Vectorizing Input and All Variables

- To help develop an efficient learning algorithm, let us *vectorize* (represent in vectors & matrices) the input and the variables involved



$$z^{[1]} = w^{[1]^T} a^{[0]} + b^{[1]}$$

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

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

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

# Neural Network Architectures

Even for a basic Neural Network, there are many design decisions to make:

1. # of hidden layers (depth)
2. # of units per hidden layer (width)
3. Type of activation function (nonlinearity)
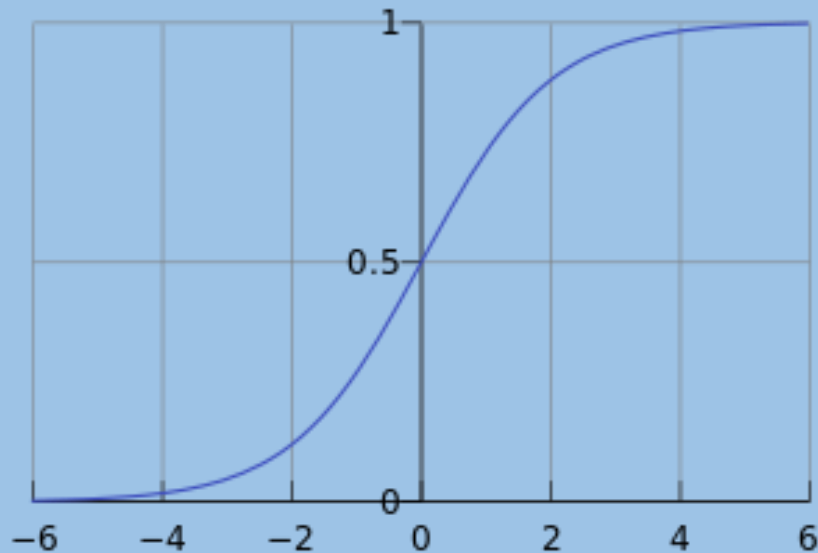4. Form of objective function

# Number of neurons

- Many neurons:
  - Higher accuracy
  - Slower
  - Risk of over-fitting
    - Memorizing, rather than understanding
    - The network will be useless with new problems.

- Few neurons:
  - Lower accuracy
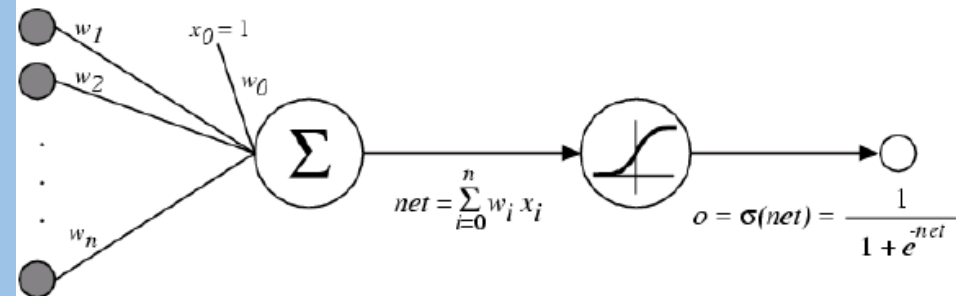  - Inability to learn at all

- Optimal number.

# Activation Functions
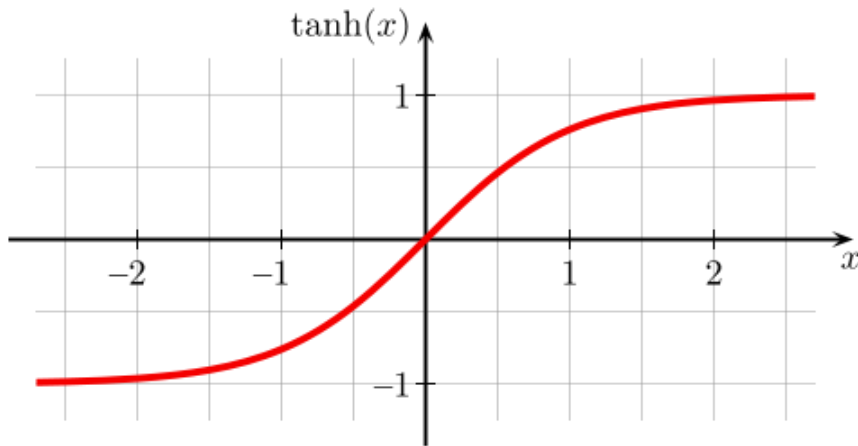
$$\text{logistic}(u) \circ \frac{1}{1+e^{-u}}$$

So far, we've assumed that the activation function (nonlinearity) is always the sigmoid function…



$$net = \sum_{i=0}^{n} w_i\, x_i$$

$$o = \sigma(net) = \frac{1}{1+e^{-net}}$$

64

# Activation Functions

- A new change: modifying the nonlinearity
  - The logistic is not widely used in modern ANNs



Alternate 1:
tanh

Like logistic function but shifted to range [-1, +1]

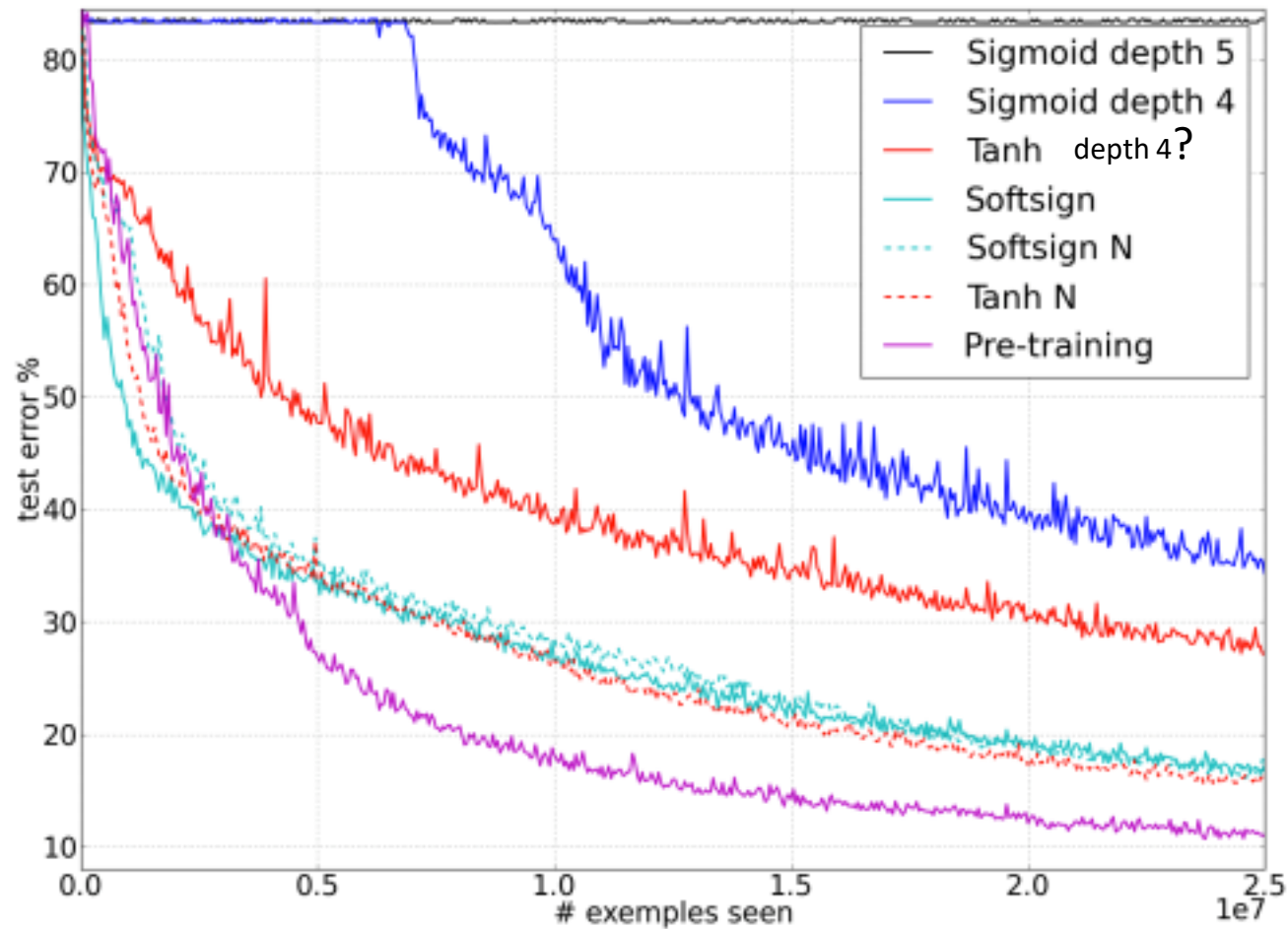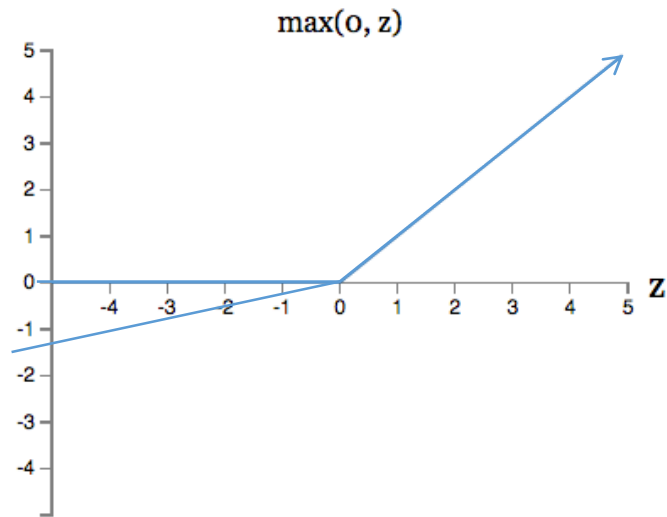# Understanding the difficulty of training deep feedforward neural networks

Figure from Glorot & Bentio (2010)

# Activation Functions

- A new change: modifying the nonlinearity
  - reLU often used in vision tasks

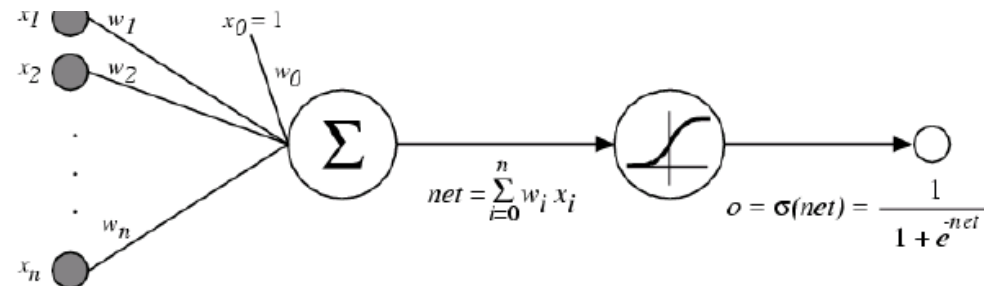$$\max(0, z)$$

$$\max(0, w \cdot x + b).$$

Alternate 2: rectified linear unit

Linear with a cutoff at zero

(Implementation: clip the gradient when you pass zero)

$$net = \sum_{i=0}^{n} w_i x_i$$

$$o = \sigma(net) = \frac{1}{1 + e^{-net}}$$

# Activation Functions

- A new change: modifying the nonlinearity
  - reLU often used in vision tasks
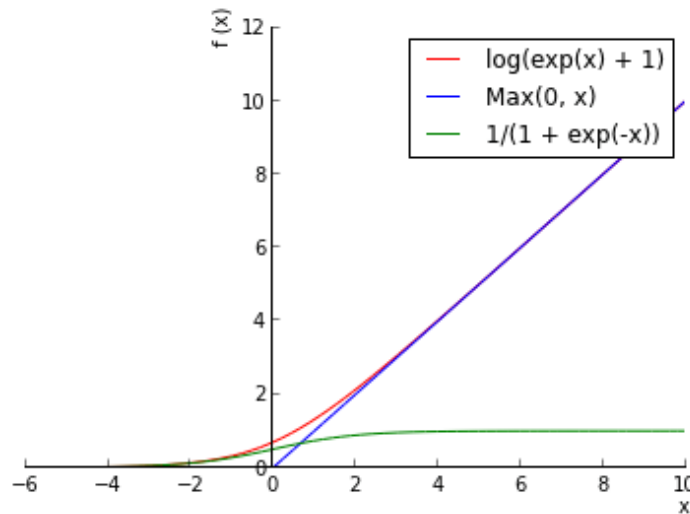


Alternate 2: rectified linear unit

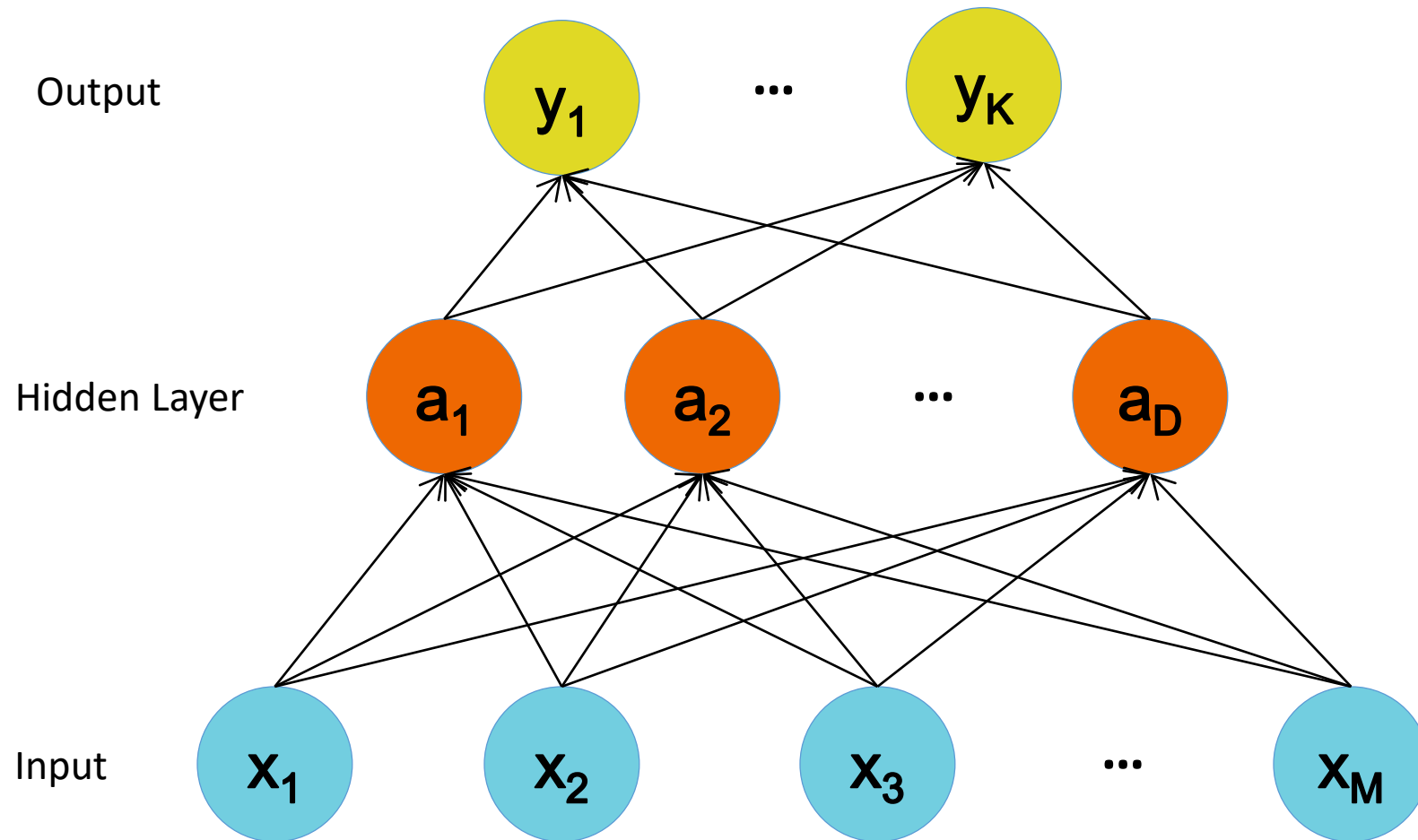Soft version: log(exp(x)+1)

Doesn't saturate (at one end)
Sparsifies outputs
Helps with vanishing gradient
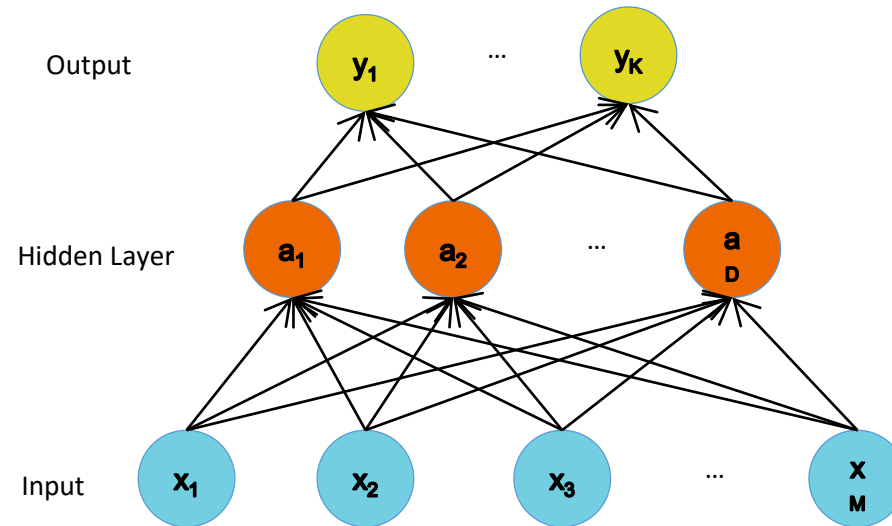
# Objective Functions for NNs

- Regression:
  - Use the same objective as Linear Regression
  - Quadratic loss (i.e. mean squared error)

- Classification:
  - Use the same objective as Logistic Regression
  - Cross-entropy (i.e. negative log likelihood)
  - This requires probabilities, so we add an additional "softmax" layer at the end of our network

# Multi-Class Output



Output

$y_1$ ... $y_K$

Hidden Layer

$a_1$ $a_2$ ... $a_D$

Input

$x_1$ $x_2$ $x_3$ ... $x_M$

# Multi-Class Output

$$y_k = \frac{\exp(b_k)}{\sum_{l=1}^{K} \exp(b_l)}$$

Softmax

# The Flow of Computations in Neural Networks

- The flow of computations in a neural network goes in two ways:
  1. Left-to-right: This is referred to as *forward propagation,* which results in computing the output of the network
  2. Right-to-left: This is referred to as *back propagation*, which results in computing the gradients (or derivatives) of the parameters in the network

- The intuition behind this 2-way flow of computations can be explained through the concept of "computation graphs"
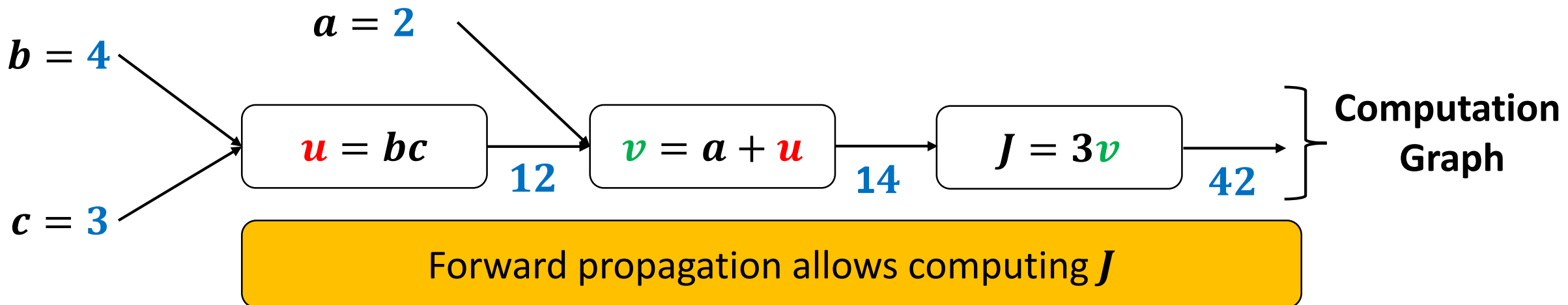
# Backpropagation

1. A set of examples for training the network is assembled. Each case consists of a problem statement (which represents the input into the network) and the corresponding solution (which represents the desired output from the network).

2. The input data is entered into the network via the input layer.

3. Each neuron in the network processes the input data with the resultant values steadily "percolating" through the network, layer by layer, until a result is generated by the output layer.

4. The actual output of the network is compared to expected output for that particular input. This results in an *error value* which represents the discrepancy between given input and expected output. On the basis of this error value an of the connection weights in the network are gradually adjusted, working backwards from the output layer, through the hidden layer, and to the input layer, until the correct output is produced. Fine tuning the weights in this way has the effect of teaching the network how to produce the correct output for a particular input, i.e. the network *learns*.
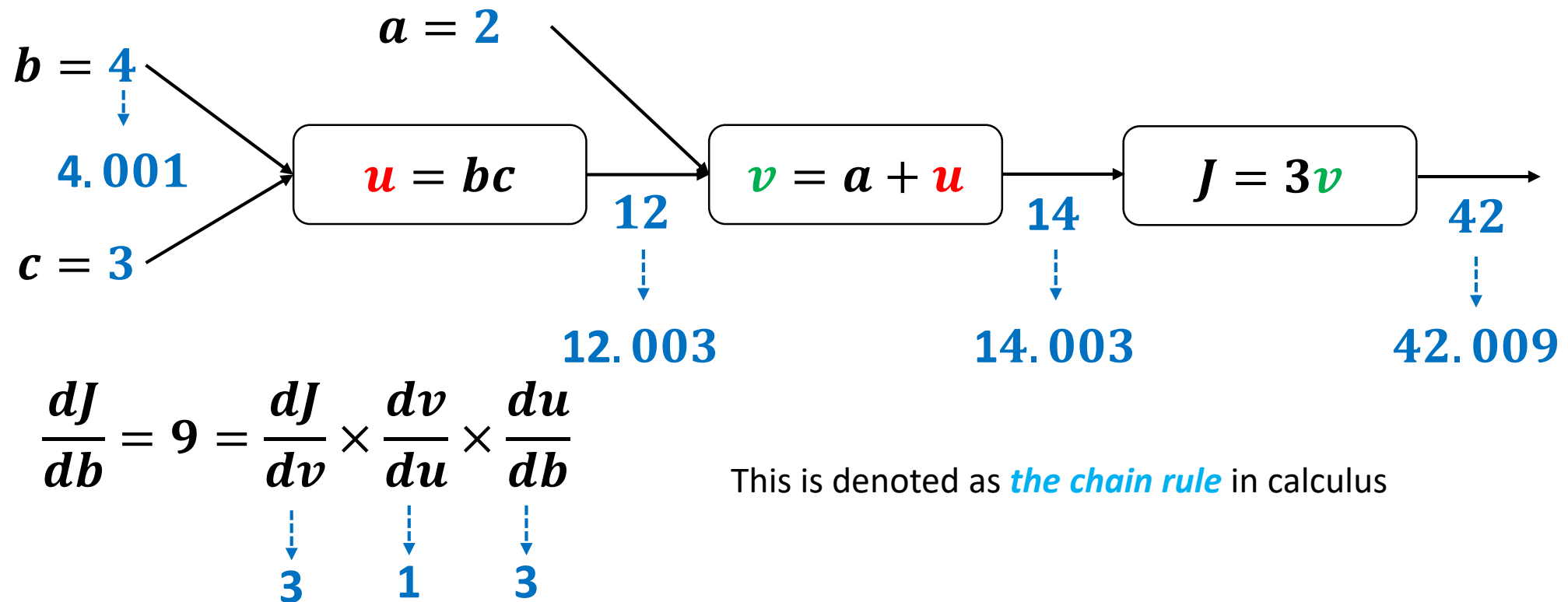
# Forward Propagation

- Let us assume we want to compute the following function $J$:

$$J(a, b, c) = 3(a + \underbrace{bc}_{u})$$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$



**Computation Graph**

Forward propagation allows computing $J$

# Backward Propagation

- Let us now compute the derivatives of the variables through the computation graph as follows:

$a = 2$

$b = 4$

$4.001$

$u = bc$

$12$

$v = a + u$

$14$

$J = 3v$

$42$

$c = 3$

$12.003$

$14.003$

$42.009$

$$\frac{dJ}{db} = 9 = \frac{dJ}{dv} \times \frac{dv}{du} \times \frac{du}{db}$$

This is denoted as *the chain rule* in calculus

$3$   $1$   $3$

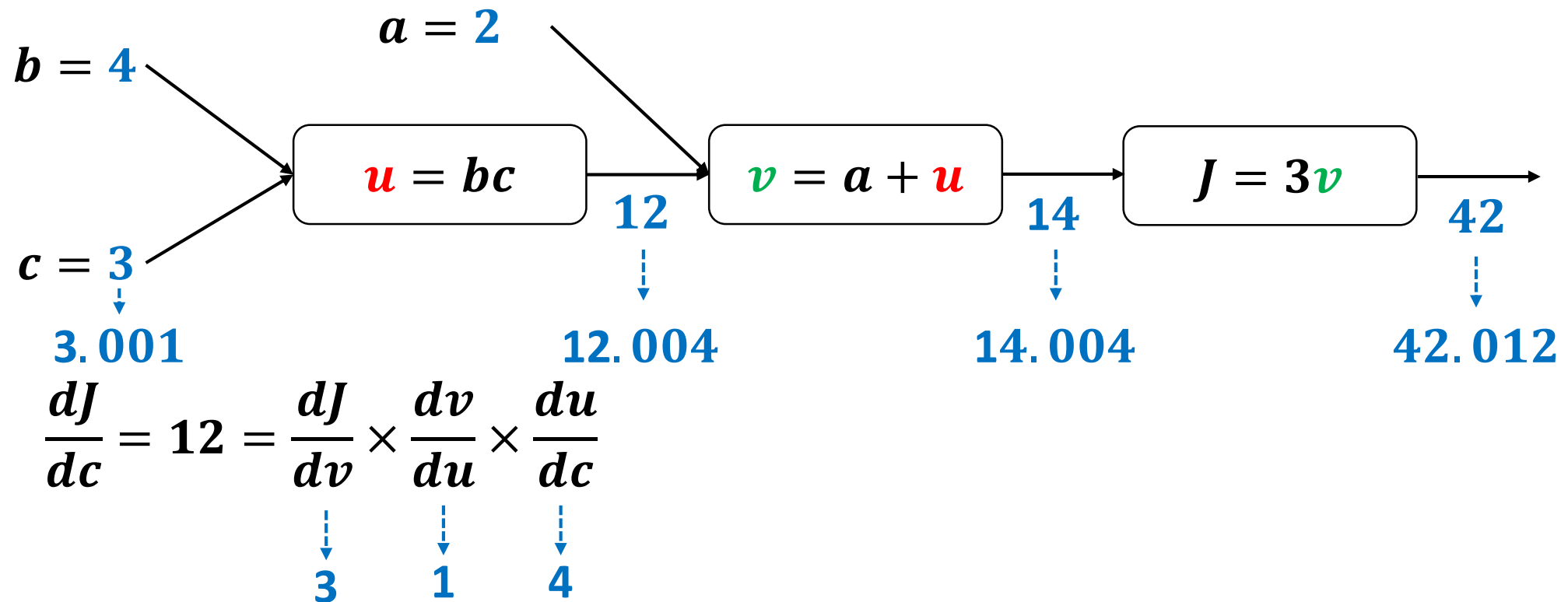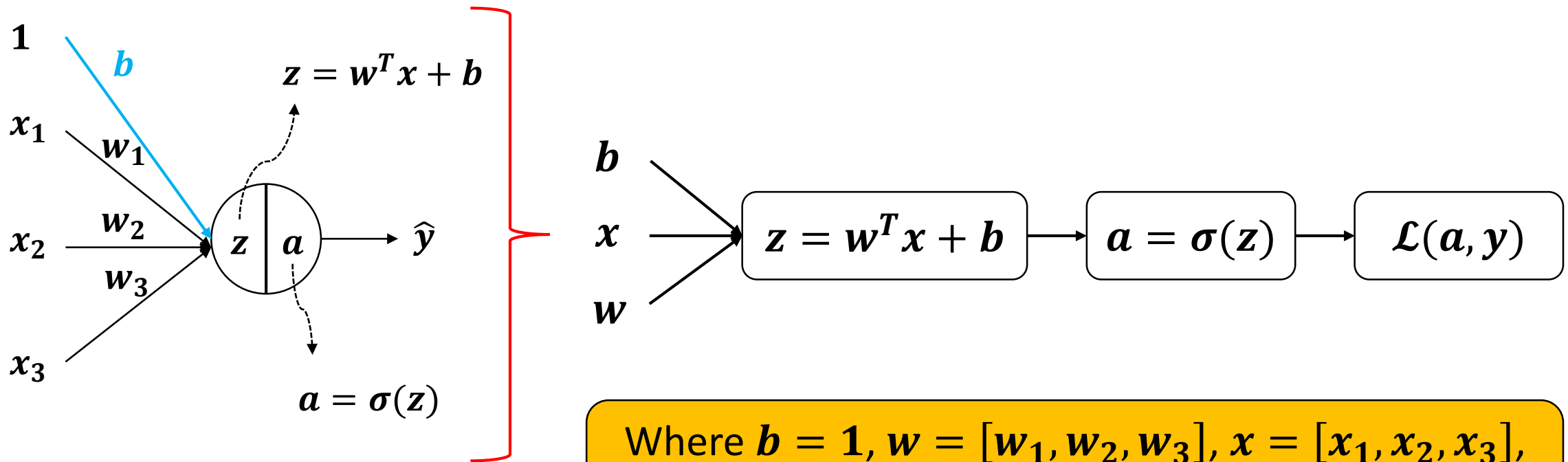# Backward Propagation

- Let us now compute the derivatives of the variables through the computation graph as follows:

$$a = 2$$

$$b = 4$$

$$c = 3$$

$$u = bc$$

$$12$$

$$v = a + u$$

$$14$$

$$J = 3v$$

$$42$$

$$3.001$$

$$12.004$$

$$14.004$$

$$42.012$$

$$\frac{dJ}{dc} = 12 = \frac{dJ}{dv} \times \frac{dv}{du} \times \frac{du}{dc}$$

$$3 \quad 1 \quad 4$$

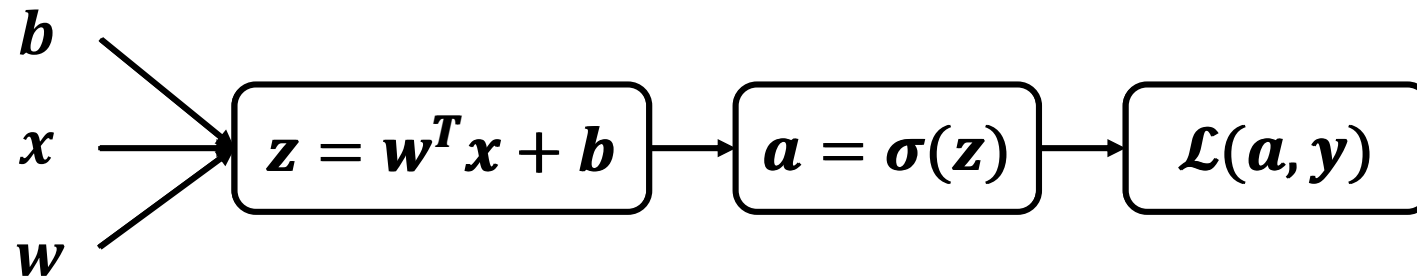# The Computation Graph of Logistic Regression

- Let us translate logistic regression (which is a **neural network** with only 1 neuron) into a computation graph



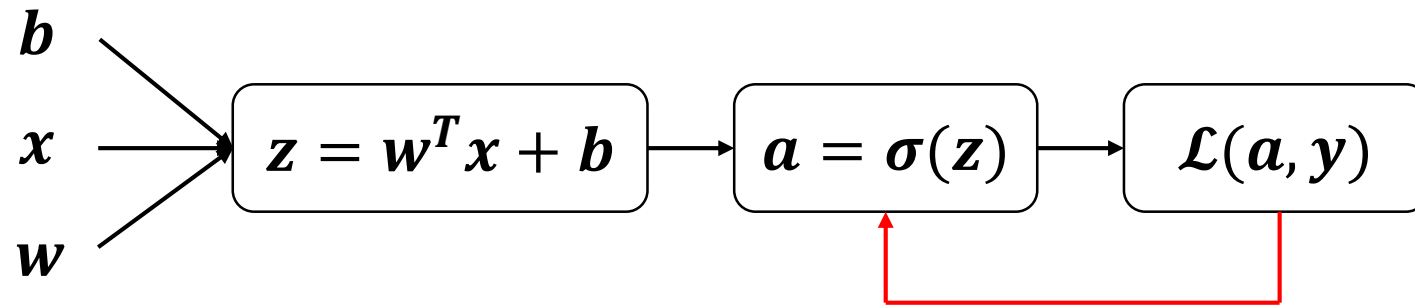Where $b = 1$, $w = [w_1, w_2, w_3]$, $x = [x_1, x_2, x_3]$, and $\mathcal{L}(a, y)$ is the cost (or *loss*) function

# Forward Propagation

- The loss function can be computed by moving from left to right

$$b$$

$$x$$

$$w$$

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

$$a = \sigma(z)$$
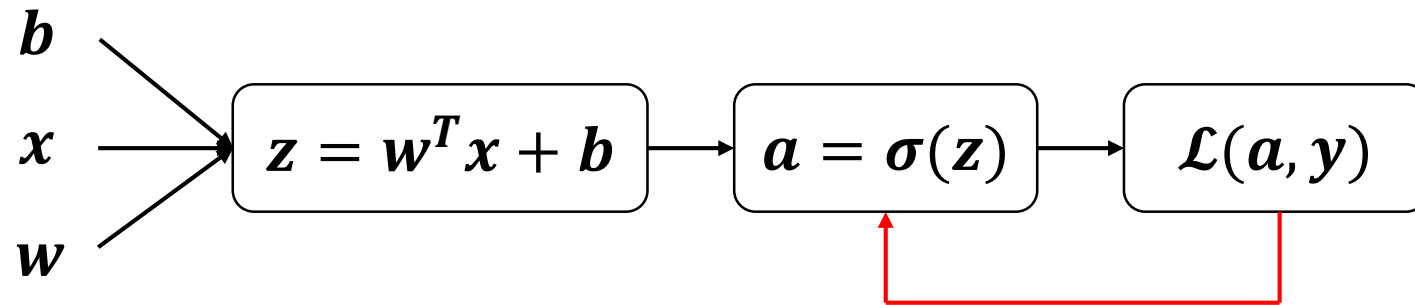
$$\mathcal{L}(a, y)$$

# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial a} = \text{Partial derivative of } \mathcal{L} \text{ with respect to } a$$
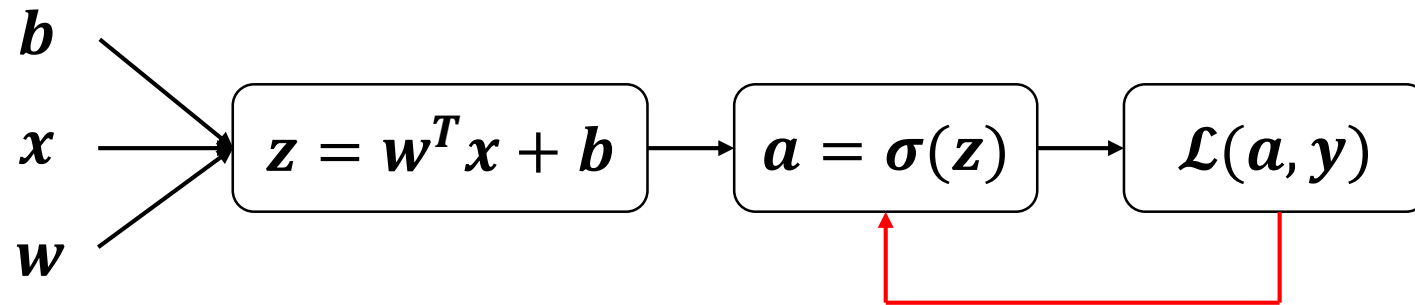
# Backward Propagation

- The derivatives can be computed by moving from right to left

$$\frac{\partial \mathcal{L}}{\partial a} = \frac{\partial}{\partial a}\left(-y\log(a) - (1-y)\log(1-a)\right)$$
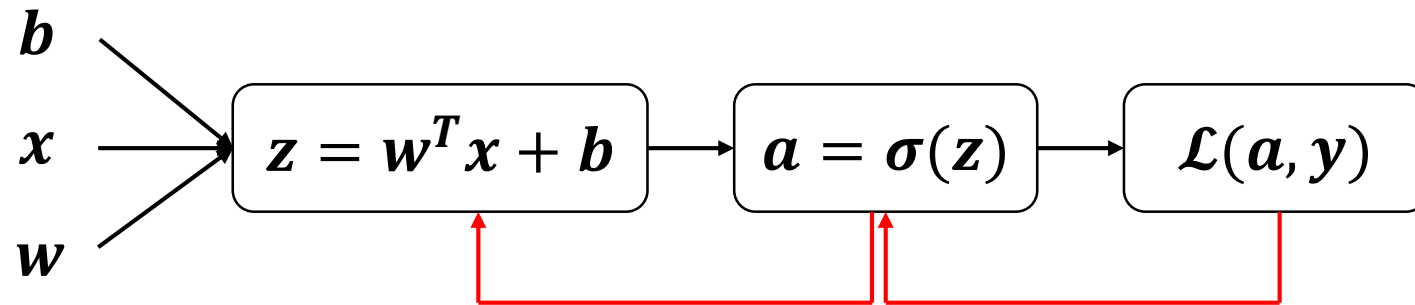
# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial a} = \frac{-y}{a} + \frac{(1-y)}{(1-a)}$$

# Backward Propagation
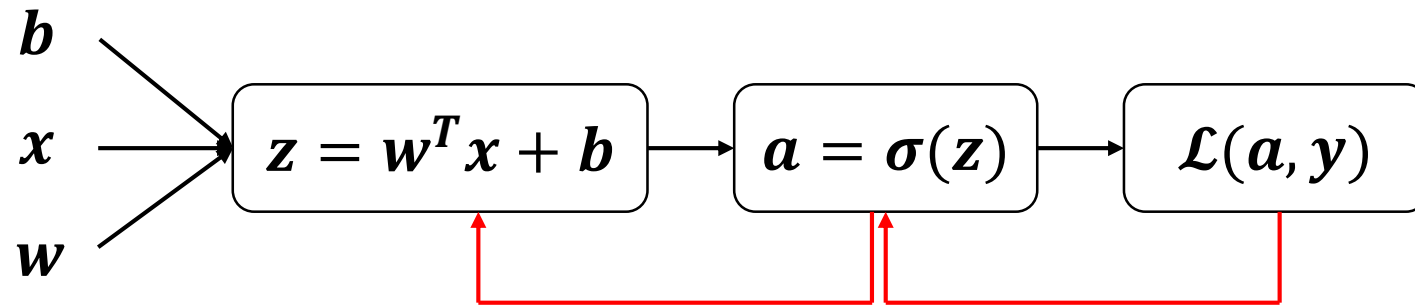
- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial z} = \text{Partial derivative of } \mathcal{L} \text{ with respect to } z$$

# Backward Propagation

- The derivatives can be computed by moving from right to left

$$z = w^T x + b \longrightarrow a = \sigma(z) \longrightarrow \mathcal{L}(a, y)$$

with inputs $b$, $x$, $w$ into the first block.

$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial a} \times \frac{\partial a}{\partial z} = \left( \frac{-y}{a} + \frac{(1-y)}{(1-a)} \right) \times \frac{\partial a}{\partial z} = \left( \frac{-y}{a} + \frac{(1-y)}{(1-a)} \right) \times a(1-a)$$

# Backward Propagation

- The derivatives can be computed by moving from right to left
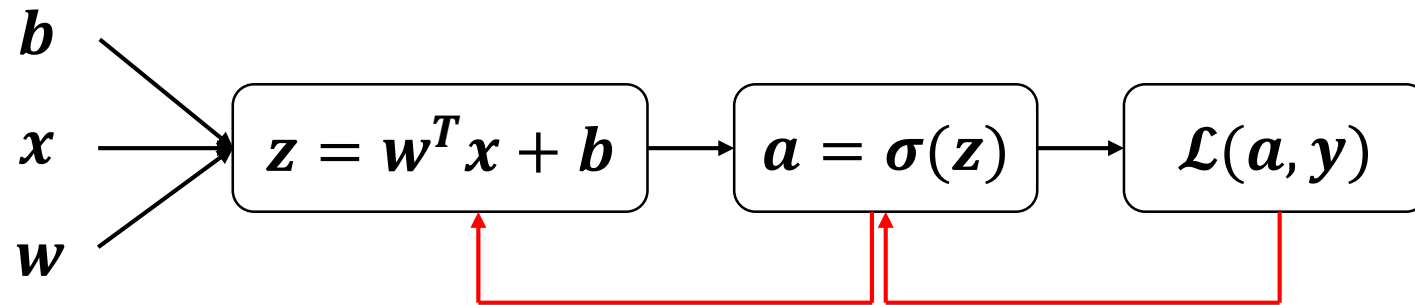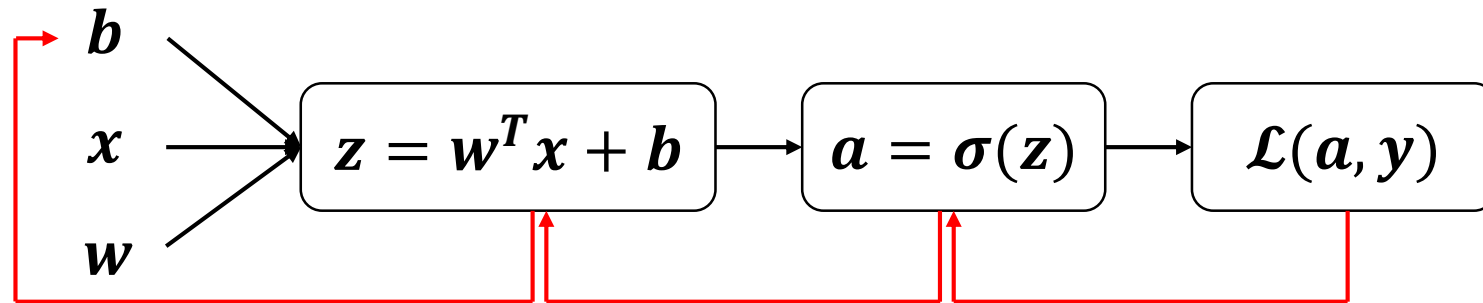


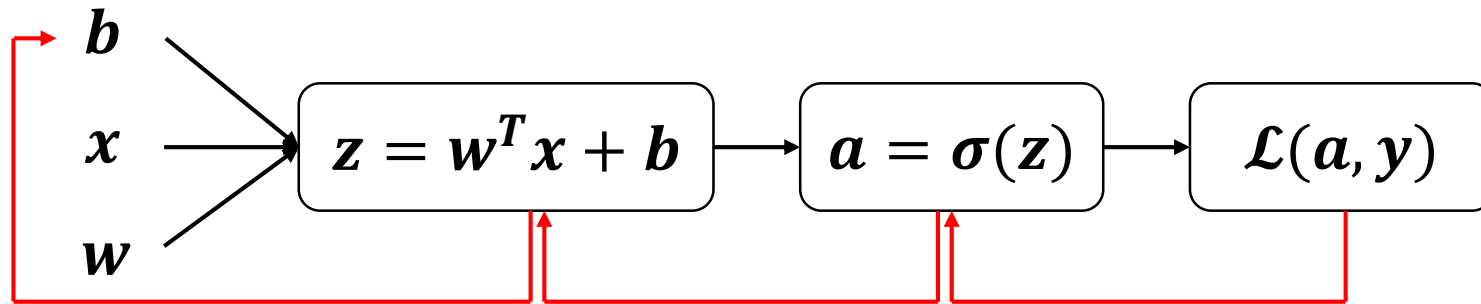$$\frac{\partial \mathcal{L}}{\partial z} = a - y$$

# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial b} = \text{Partial derivative of } \mathcal{L} \text{ with respect to } b$$
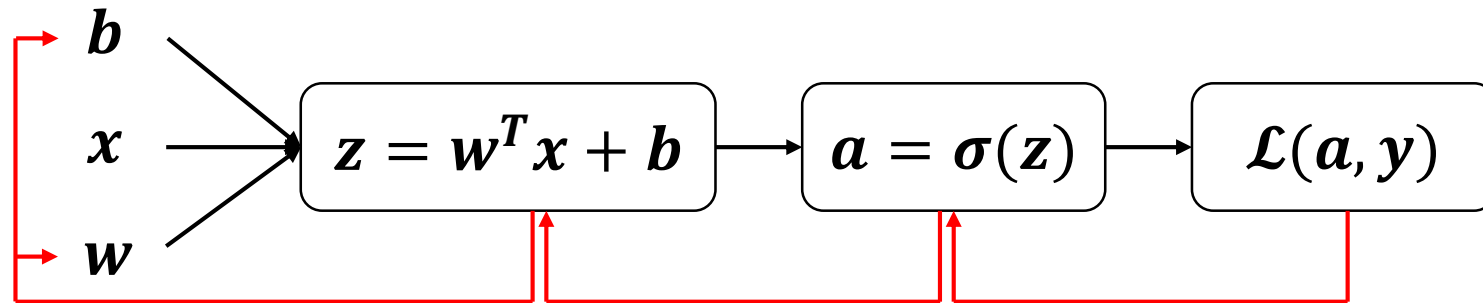
# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial b} = (a - y) \times \frac{\partial z}{\partial b} = (a - y) \times \mathbf{1} = (a - y)$$
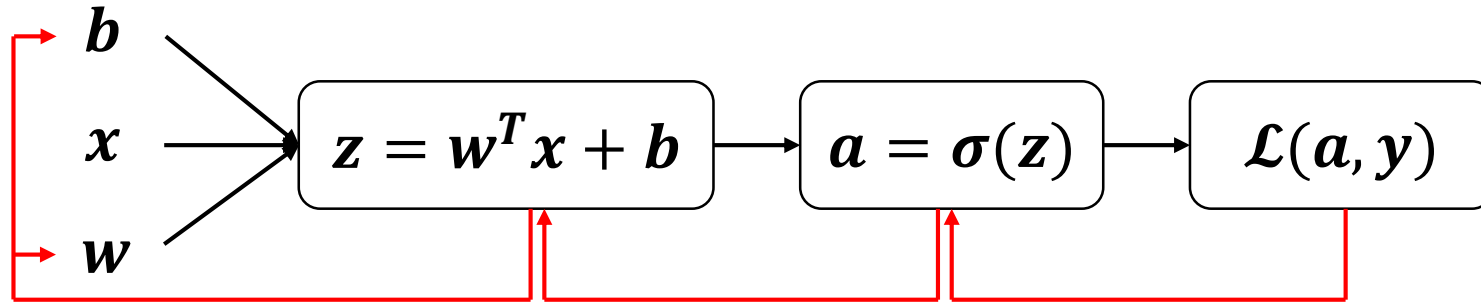
# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial w} = \text{Partial derivative of } \mathcal{L} \text{ with respect to } w$$
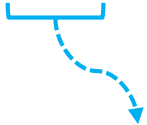
# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w} = (a - y) \times \frac{\partial z}{\partial w} = (a - y)x$$

# Backward Propagation: Summary

- Here is the summary of the gradients in logistic regression:

$$dz = \frac{\partial \mathcal{L}}{\partial z} = a - y$$

We will denote this as $dz$ for simplicity

# Backward Propagation: Summary

- Here is the summary of the gradients in logistic regression:

$$dz = \frac{\partial \mathcal{L}}{\partial z} = a - y$$

$$db = \frac{\partial \mathcal{L}}{\partial b} = a - y$$

$$dw = \frac{\partial \mathcal{L}}{\partial w} = (a - y)x$$

# Gradient Descent For Logistic Regression

- **Outline**:
  - Have a loss function $\mathcal{L}(w, b)$, where $w = [w_1, \ldots, w_m]$ and $b = w_0$
  - Start off with some guesses for $w_1, \ldots, w_m$
    - It does not really matter what values you start off with for $w_1, \ldots, w_m$, but a common choice is to set them all initially to zero
  - Repeat until convergence{

$$w_j = w_j - \alpha \frac{\partial \mathcal{L}(w, b)}{\partial w_j}$$

$$b = b - \alpha \frac{\partial \mathcal{L}(w, b)}{\partial b}$$

**Let us focus on this part**

  }

# Gradient Descent For Logistic Regression

- **Outline**:

Assuming $n$ examples

- Repeat until convergence{

$$for\ i = 1\ to\ n:$$

**Forward propagation**
$$z^{(i)} = w^T x^{(i)} + b$$
$$a^{(i)} = \sigma(z^{(i)})$$

**Backward propagation**
$$dz^{(i)} = a^{(i)} - y^{(i)}$$
$$dw = dw + dz^{(i)} x^{(i)}$$
$$db = db + dz^{(i)}$$

**Outside the loop**
$$dw = dw/n$$
$$db = db/n$$
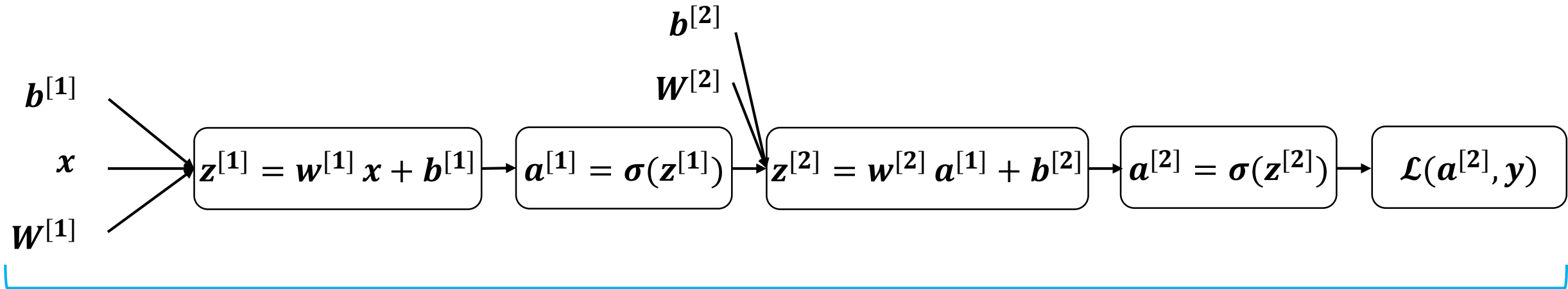$$w = w - \alpha dw$$
$$b = b - \alpha db$$
}

# The Computation Graph of A Neural Network

- we can represent any neural network in terms of a computation graph
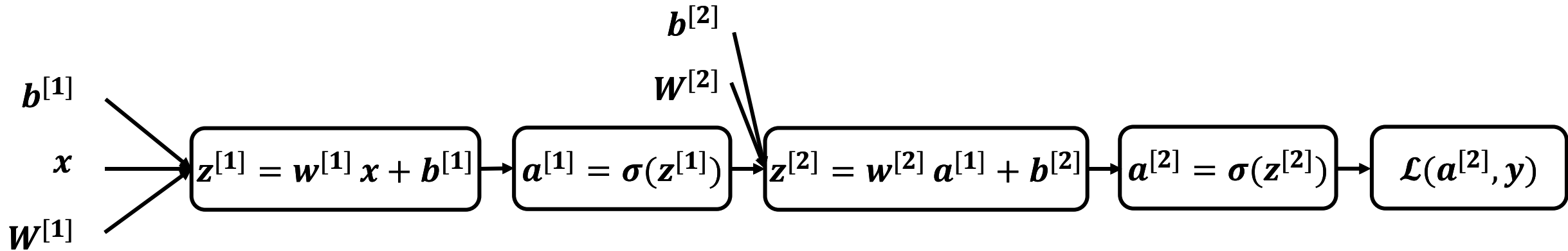
# The Computation Graph of A Neural Network

- Akin to logistic regression, we can represent any neural network in terms of a computation graph



**The corresponding computation graph**

# Forward Propagation

- The loss function can be computed by moving from left to right

$$b^{[2]}$$

$$b^{[1]}$$

$$W^{[2]}$$

$$x$$

$$W^{[1]}$$

$$z^{[1]} = w^{[1]} x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow \mathcal{L}(a^{[2]}, y)$$
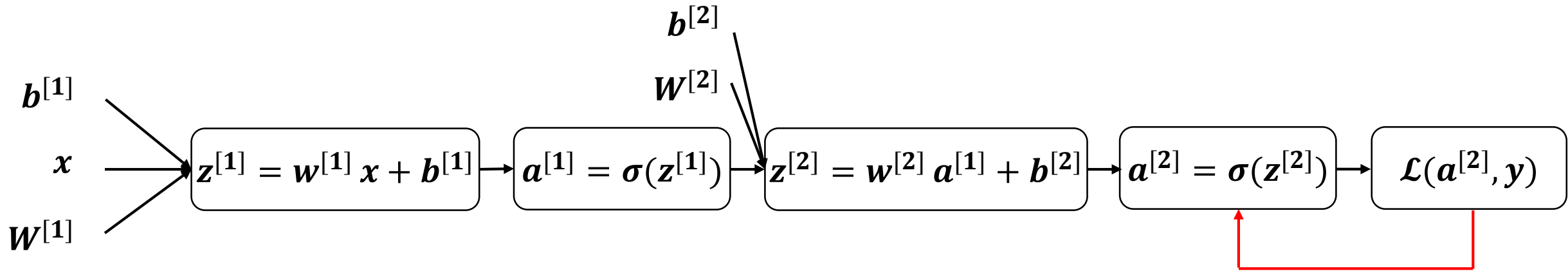
# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial a^{[2]}} = \text{Partial derivative of } \mathcal{L} \text{ with respect to } a^{[2]}$$
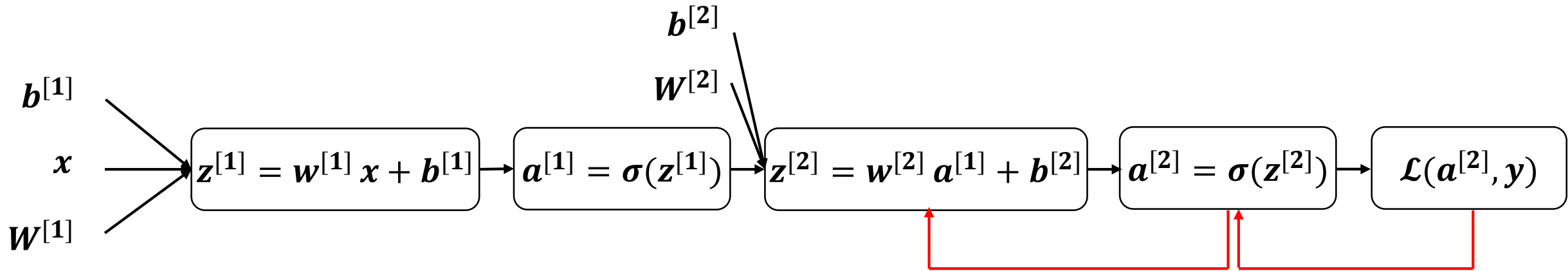
# Backward Propagation

- The derivatives can be computed by moving from right to left

$$b^{[2]}$$

$$W^{[2]}$$

$$b^{[1]}$$

$$x$$

$$W^{[1]}$$

$$z^{[1]} = w^{[1]} x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow \mathcal{L}(a^{[2]}, y)$$

$$\frac{\partial \mathcal{L}}{\partial a^{[2]}} = \frac{-y}{a^{[2]}} + \frac{(1-y)}{(1-a^{[2]})}$$
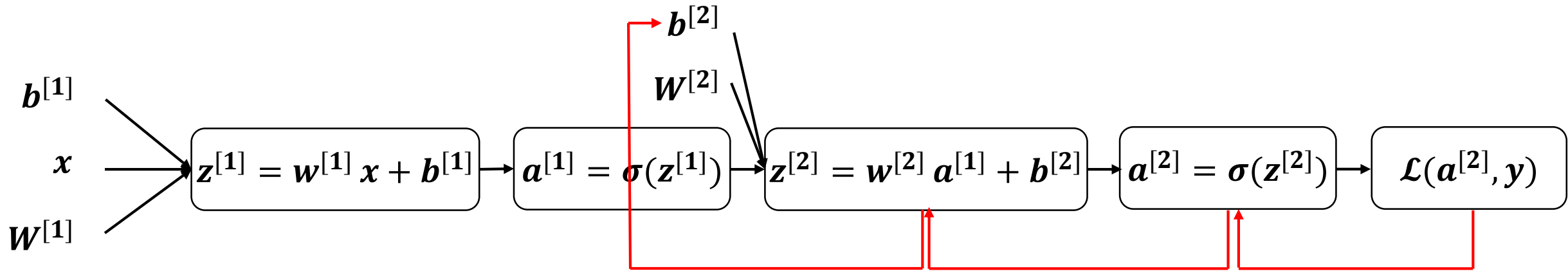
# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial z^{[2]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} = a^{[2]} - y$$
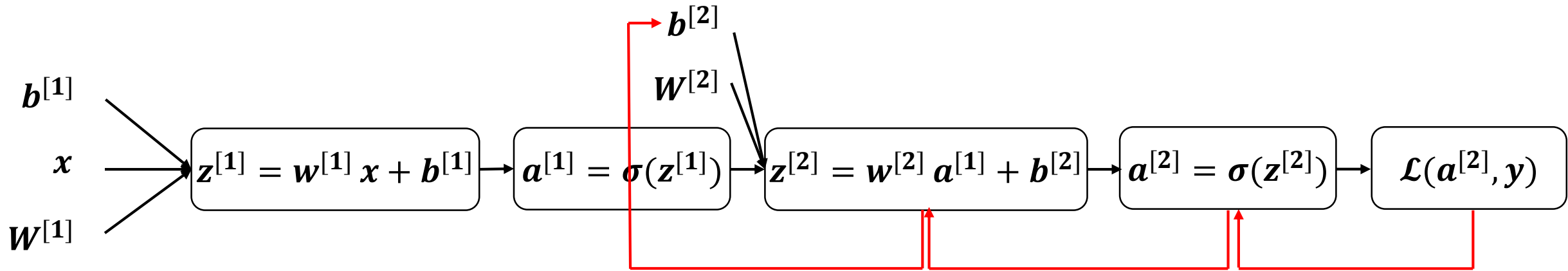
# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial b^{[2]}} = \text{Partial derivative of } \mathcal{L} \text{ with respect to } b^{[2]}$$
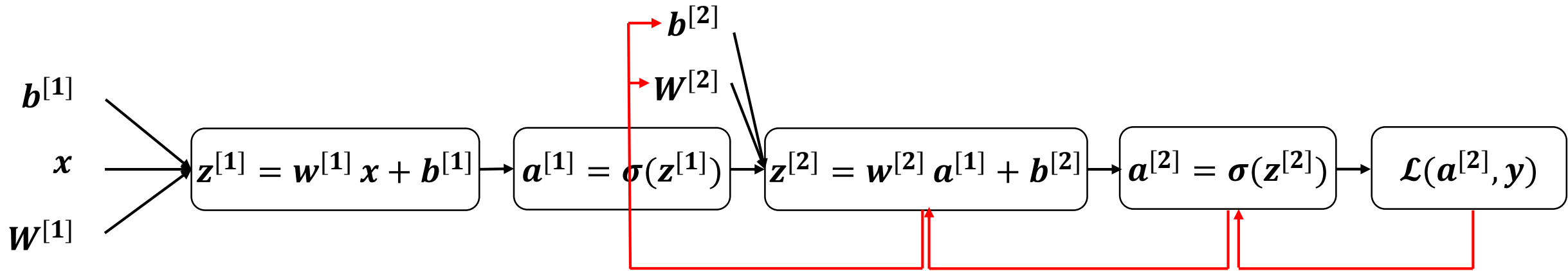
# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial b^{[2]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial b^{[2]}} = a^{[2]} - y$$
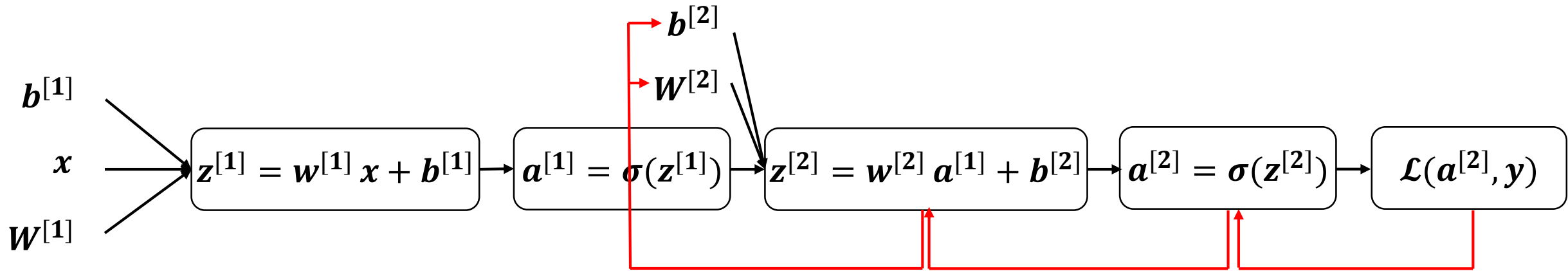
# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial W^{[2]}} = \text{Partial derivative of } \mathcal{L} \text{ with respect to } W^{[2]}$$

# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial W^{[2]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial W^{[2]}} = (a^{[2]} - y)a^{[1]T}$$
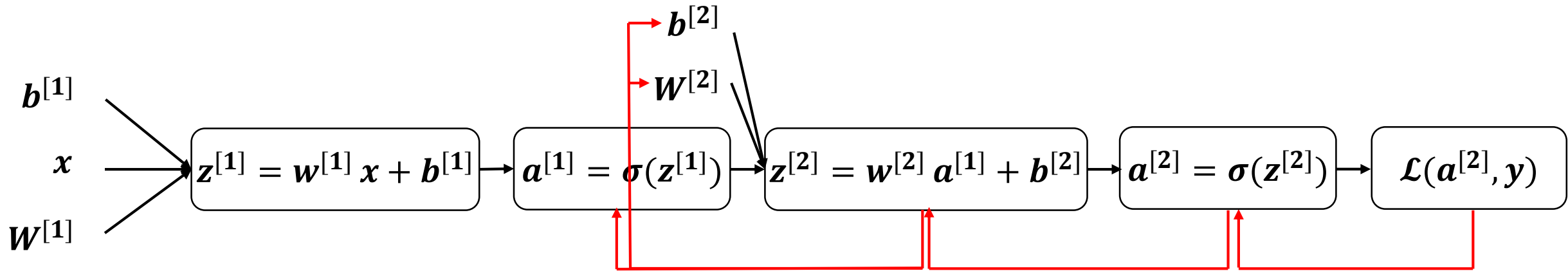
# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial a^{[1]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial a^{[1]}} = (a^{[2]} - y)w^{[2]T}$$

# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial z^{[1]}} = \text{ Partial derivative of } \mathcal{L} \text{ with respect to } z^{[1]}$$
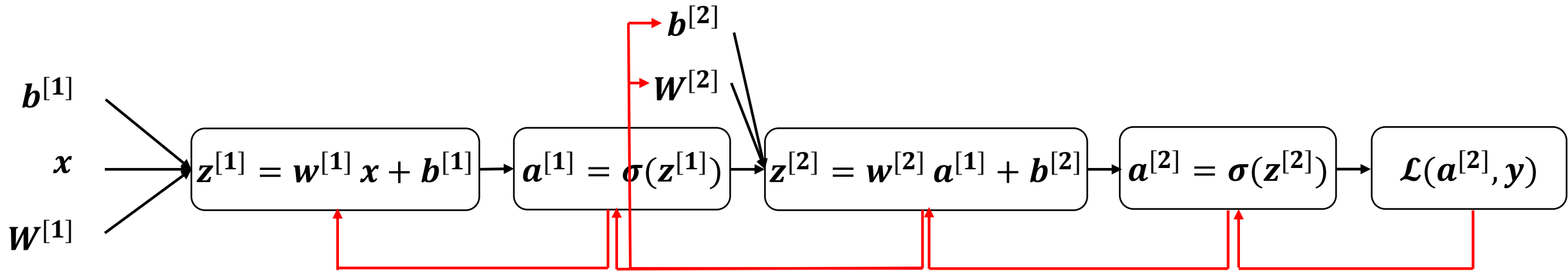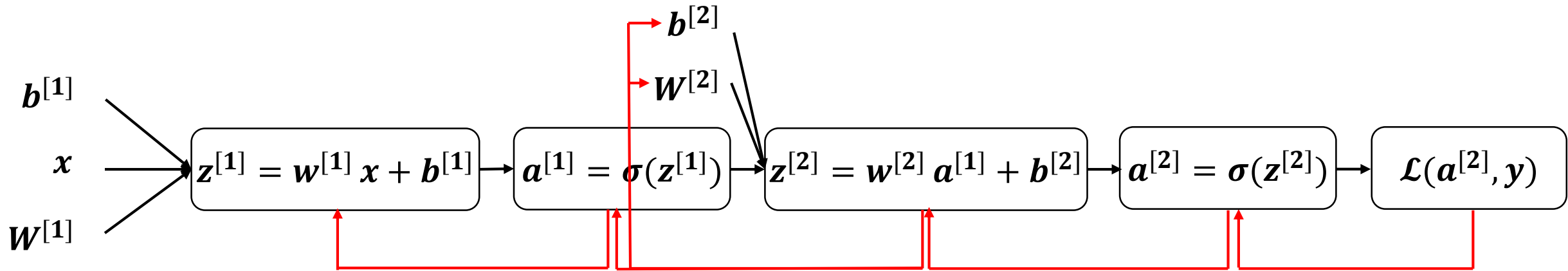
# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial z^{[1]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial a^{[1]}} \times \frac{\partial a^{[1]}}{\partial z^{[1]}} = (a^{[2]} - y)w^{[2]T} * a^{[1]}(1 - a^{[1]})$$

Element-wise product

# Backward Propagation

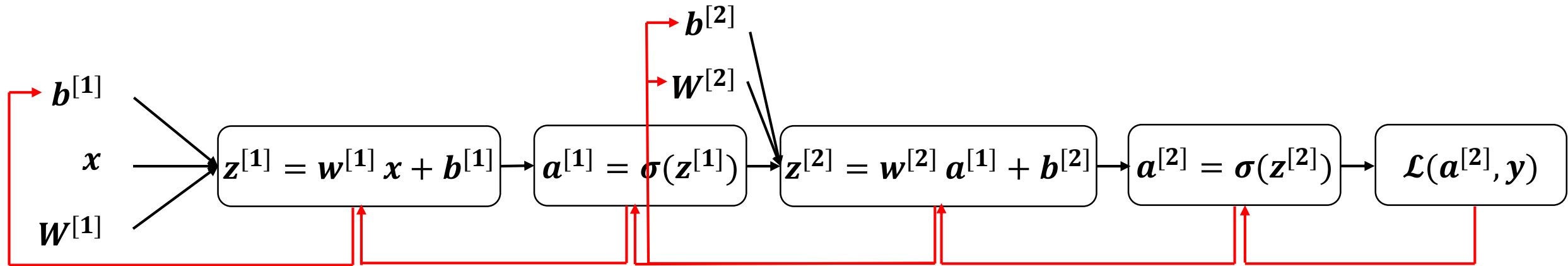- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial b^{[1]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial a^{[1]}} \times \frac{\partial a^{[1]}}{\partial z^{[1]}} \times \frac{\partial z^{[1]}}{\partial b^{[1]}}$$

# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial b^{[1]}} = \; (a^{[2]} - y)w^{[2]T} * a^{[1]}(1 - a^{[1]})$$
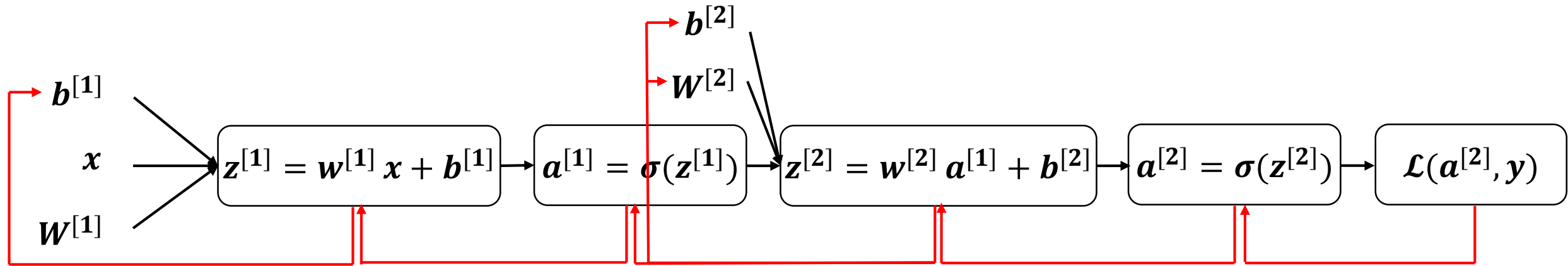
# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial W^{[1]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial a^{[1]}} \times \frac{\partial a^{[1]}}{\partial z^{[1]}} \times \frac{\partial z^{[1]}}{\partial W^{[1]}}$$
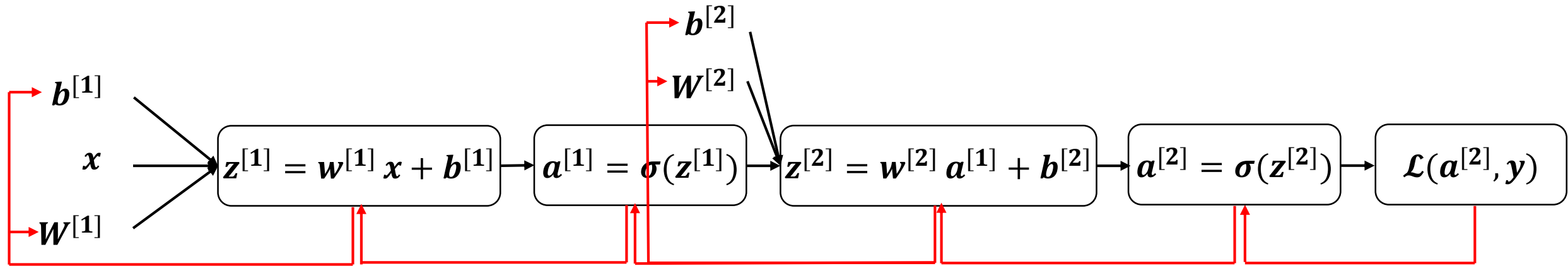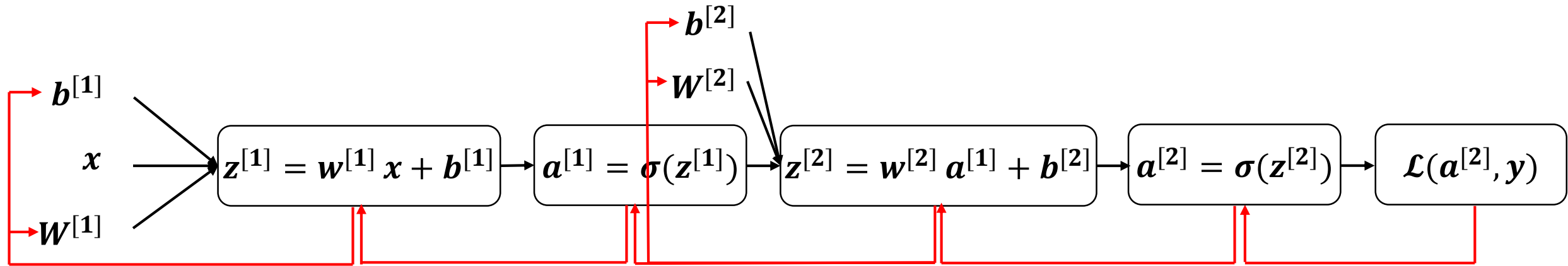
# Backward Propagation

- The derivatives can be computed by moving from right to left



$$\frac{\partial \mathcal{L}}{\partial W^{[1]}} = \left( (a^{[2]} - y) w^{[2]T} * a^{[1]} (1 - a^{[1]}) \right) x^T$$

# Backward Propagation: Summary

- Here is the summary of the gradients in our given neural network:

$$dz^{[2]} = \frac{\partial \mathcal{L}}{\partial z^{[2]}} = a^{[2]} - y \qquad\qquad dz^{[1]} = \frac{\partial \mathcal{L}}{\partial z^{[1]}} = dz^{[2]} w^{[2]T} * a^{[1]}(1 - a^{[1]})$$

$$db^{[2]} = \frac{\partial \mathcal{L}}{\partial b^{[2]}} = a^{[2]} - y \qquad\qquad db^{[1]} = \frac{\partial \mathcal{L}}{\partial b^{[1]}} = dz^{[2]} w^{[2]T} * a^{[1]}(1 - a^{[1]})$$

$$dW^{[2]} = \frac{\partial \mathcal{L}}{\partial W^{[2]}} = (a^{[2]} - y)a^{[1]T} \qquad dW^{[1]} = \frac{\partial \mathcal{L}}{\partial W^{[1]}} = \left( dz^{[2]} w^{[2]T} * a^{[1]}(1 - a^{[1]}) \right) x^{T}$$

# Perceptron

- **Architecture du réseau**

- Nous importons les classes Sequential et Dense pour définir notre modèle et son architecture.

```
#keras
from keras.models import Sequential
from keras.layers import Dense
```

- La classe Sequential est une structure, initialement vide, qui permet de définir un empilement de couches de neurones (https://keras.io/getting-started/sequential-model-guide/)

# Perceptron multicouche

- Dans cette section, nous avons un perceptron multicouche. Nous créons toujours une structure Sequential, dans lequel nous ajoutons successivement deux objets Dense : le premier fait la jonction entre la couche d'entrée (d'où l'option input_dim indiquant le nombre de variables prédictives) et la couche cachée qui comporte (units = 3) neurones ; le second entre cette couche cachée et la sortie à un seul neurone (units = 1). Nous avons une fonction d'activation sigmoïde dans les deux cas.

# MLP en utilisant Keras

```
#modélisation
modelMc = Sequential()
modelMc.add(Dense(units=3,input_dim=2,activation="sigmoid"))
modelMc.add(Dense(units=1,activation="sigmoid"))
```
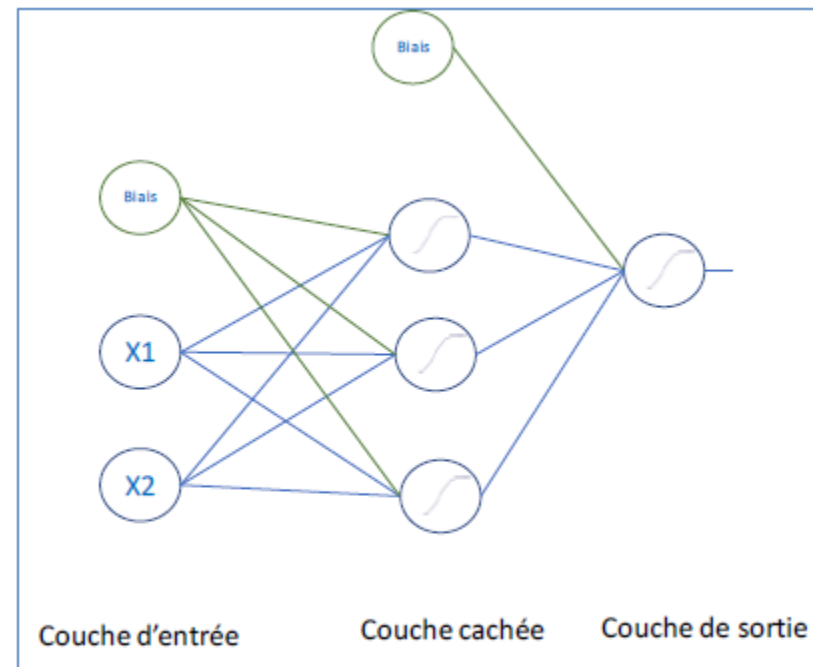


Figure 4 - Perceptron multicouche - Structure

```python
#compilation - algorithme d'apprentissage
modelMc.compile(loss="binary_crossentropy",optimizer="adam",metrics=["accuracy"])

#apprentissage
modelMc.fit(XTrain,yTrain,epochs=150,batch_size=10)

#poids synaptiques
print(modelMc.get_weights())
```

```python
#score
score = modelMc.evaluate(XTest,yTest)
print(score)
```
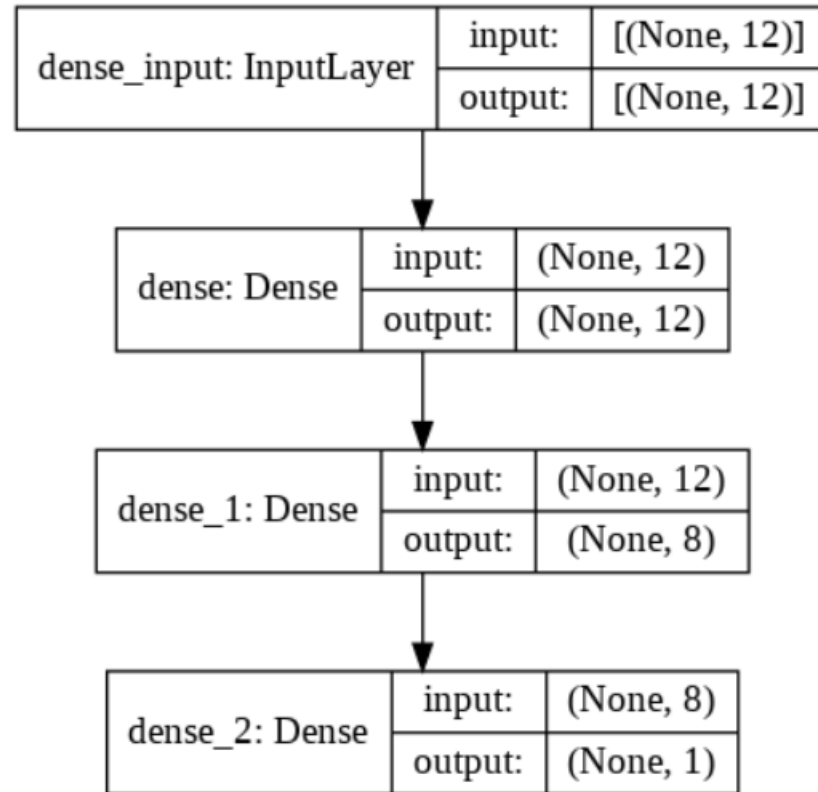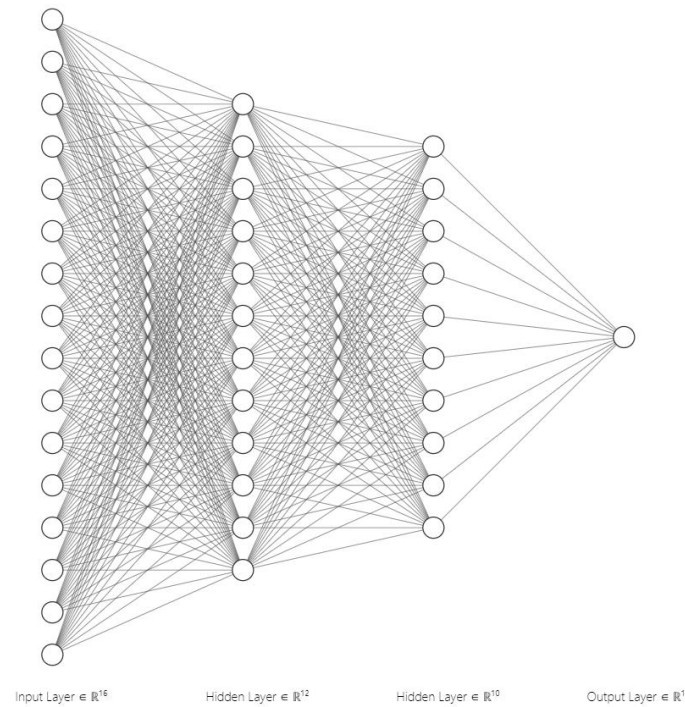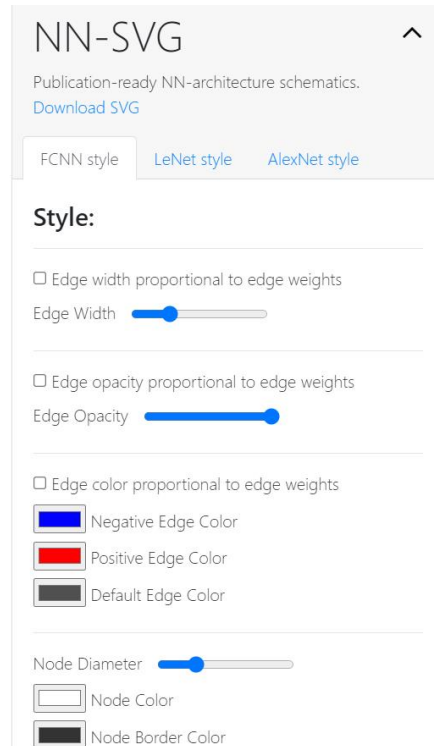
# Here two hidden layers..

g

```python
# Initializing the ANN
ann = tf.keras.models.Sequential()
# Add the input layer and first hidden layer
ann.add(tf.keras.layers.Dense(units=12, activation='relu', input_shape=X_t
# Add the second hidden layer
ann.add(tf.keras.layers.Dense(units=8, activation='relu'))
# Add the output layer
ann.add(tf.keras.layers.Dense(units=1, activation='sigmoid'))
```

```python
from tensorflow.keras.utils import plot_model
plot_model(ann,
           to_file="model.png",
           show_shapes=True,
           show_layer_names=True,
          )
```

**Output**

| dense_input: InputLayer | input: | [(None, 12)] |
| | output: | [(None, 12)] |

| dense: Dense | input: | (None, 12) |
| | output: | (None, 12) |

| dense_1: Dense | input: | (None, 12) |
| | output: | (None, 8) |

| dense_2: Dense | input: | (None, 8) |
| | output: | (None, 1) |

# https://alexlenail.me/NN-SVG/