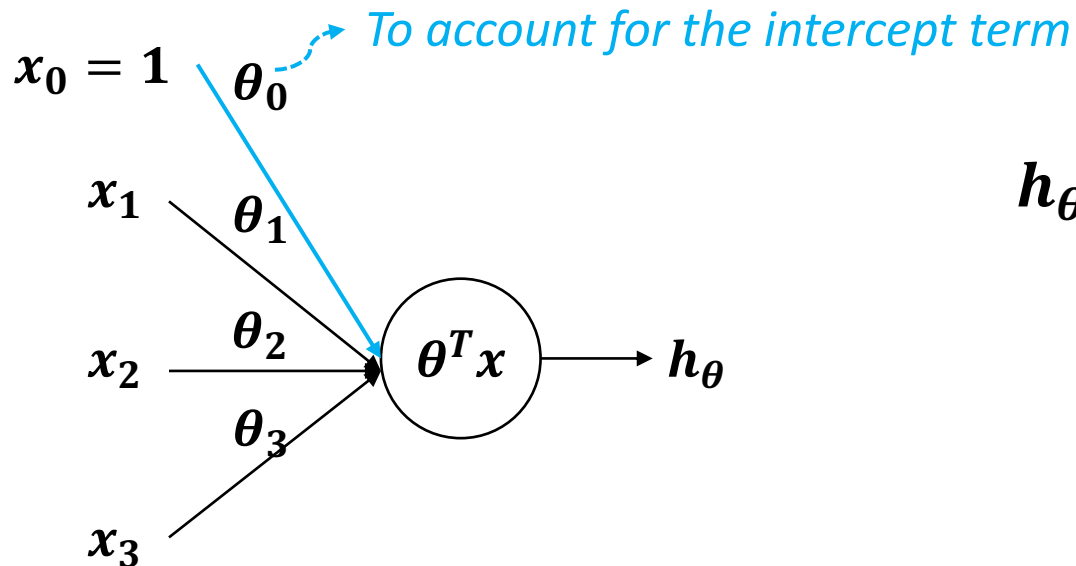


# Basics of Neural Networks

# Linear Regression

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



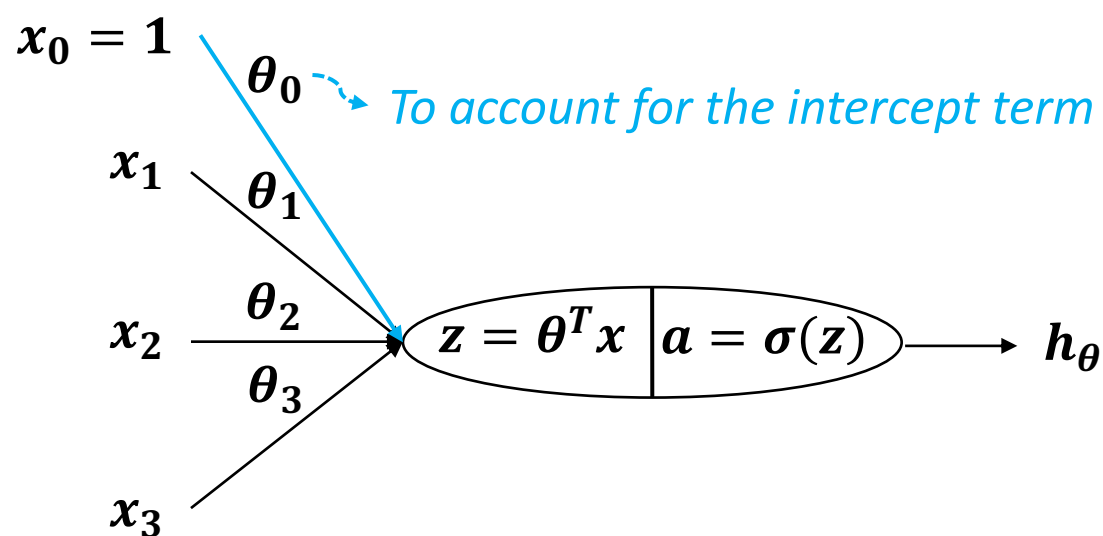
$$h_{\theta} = \theta^T x = [\theta_0 \ \theta_1 \ \theta_2 \ \theta_3] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

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

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

# Logistic Regression

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



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

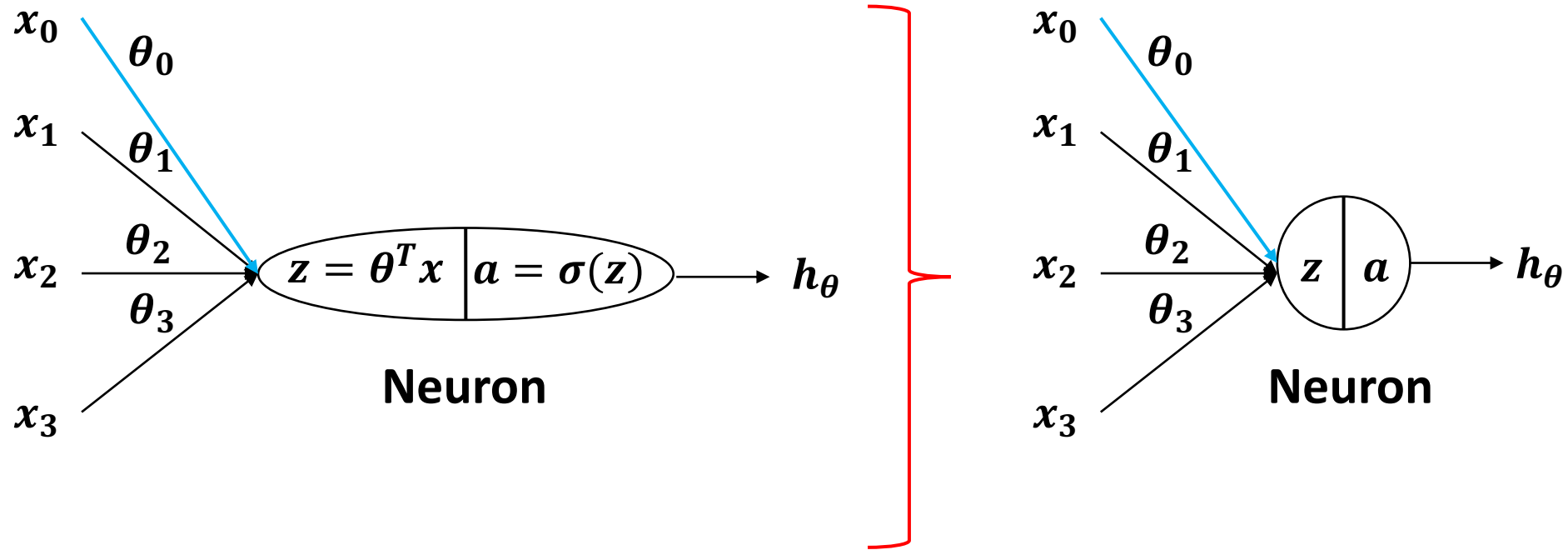
$$= \sigma(\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3)$$

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

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

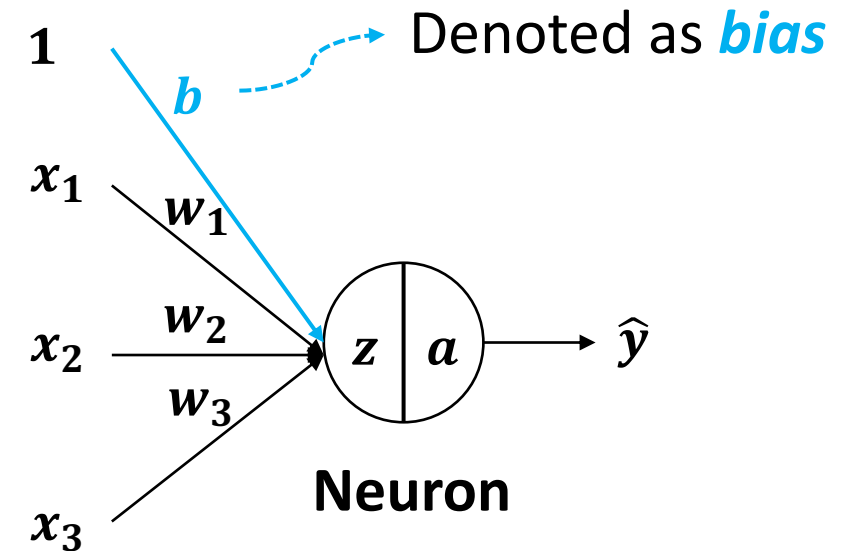
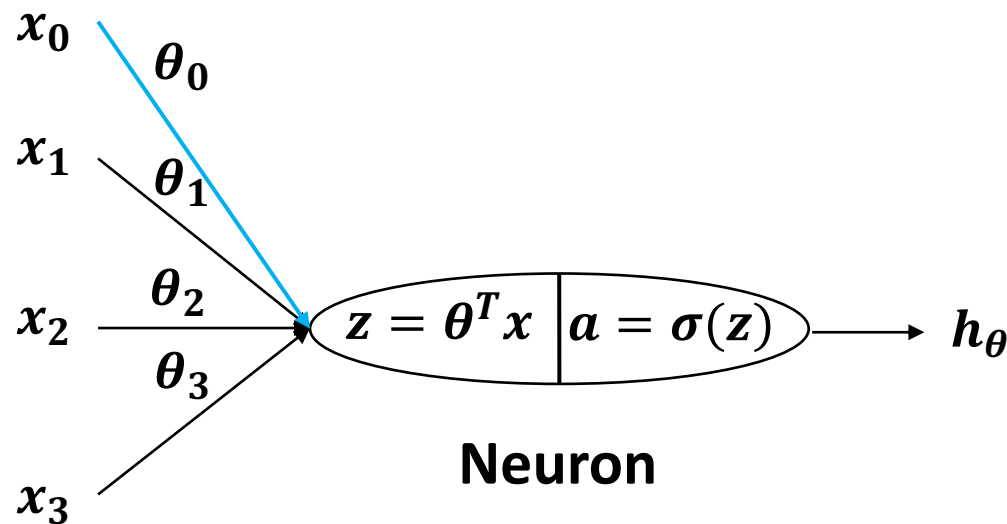
# Towards Neural Networks

Logistic regression can be considered as **neural network** with only 1 neuron



# Towards Neural Networks

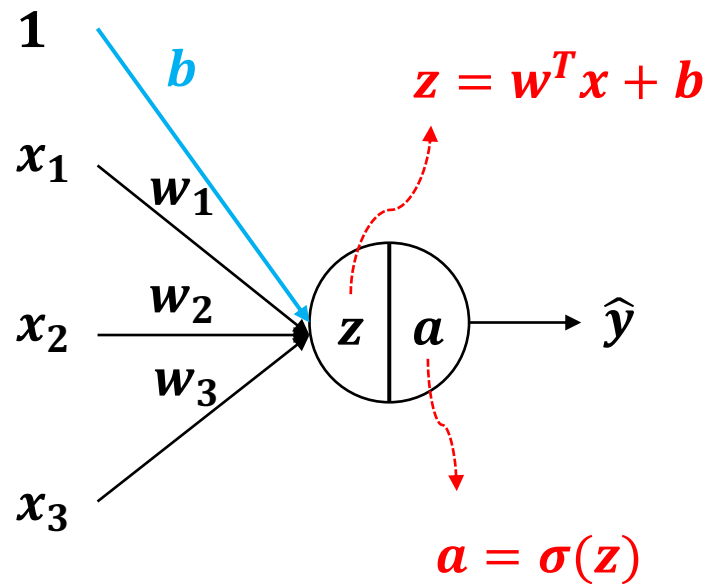
Logistic regression can be considered as **neural network** with only 1 neuron



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

# Towards Neural Networks

Logistic regression can be considered as **neural network** with only 1 neuron



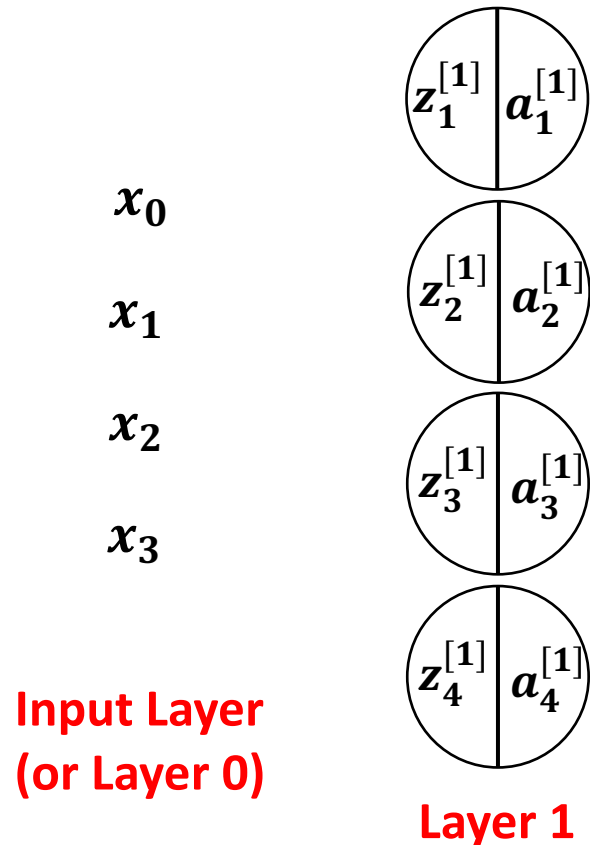
$$\hat{y} = a = \sigma(z) = \sigma(w^T x + b) = \sigma([w_1 \ w_2 \ w_3] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + b)$$

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

$$= \frac{1}{1 + e^{-(w_1 x_1 + w_2 x_2 + w_3 x_3 + b)}}$$

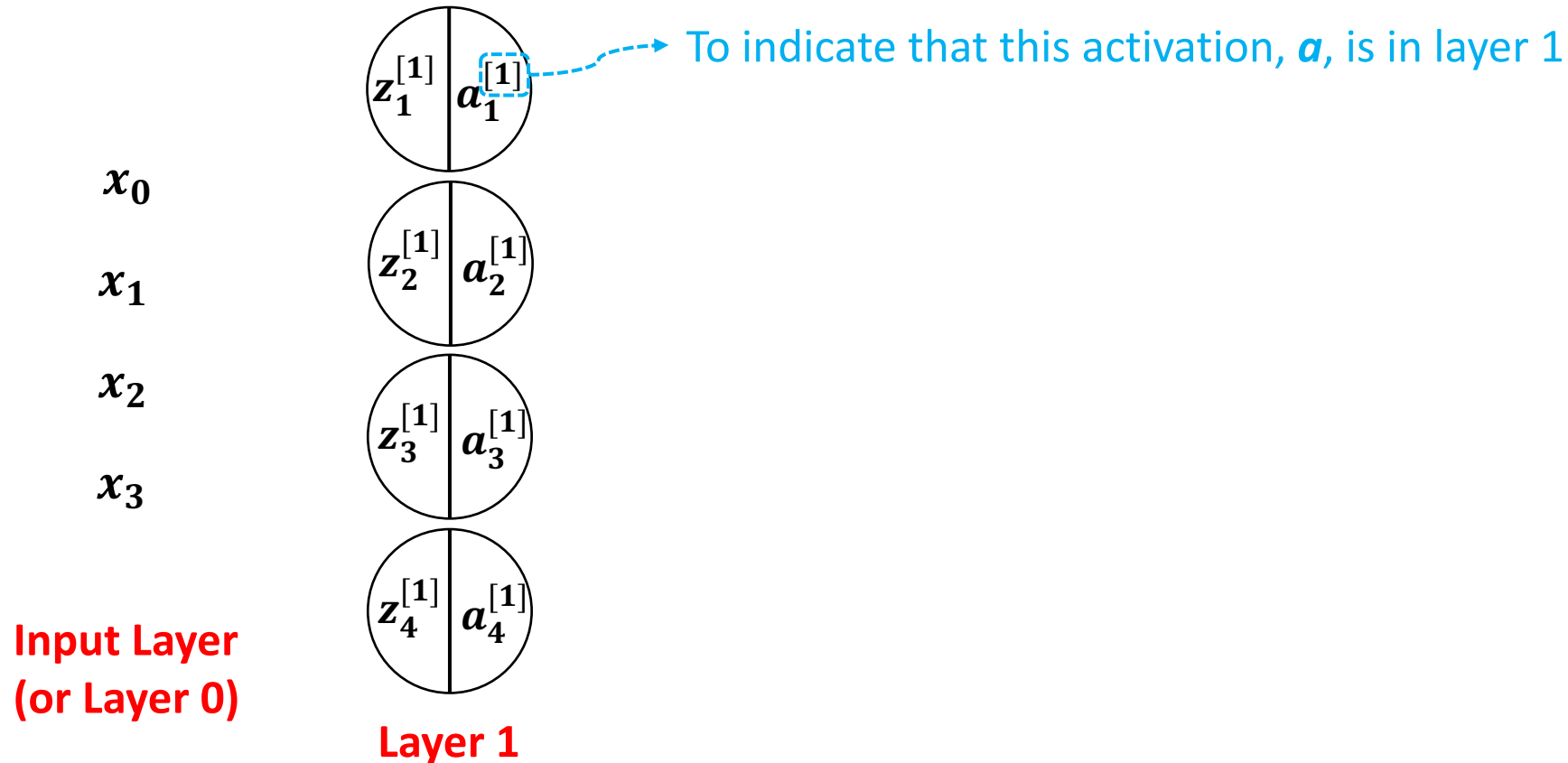
# Neural Networks

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



# Neural Networks

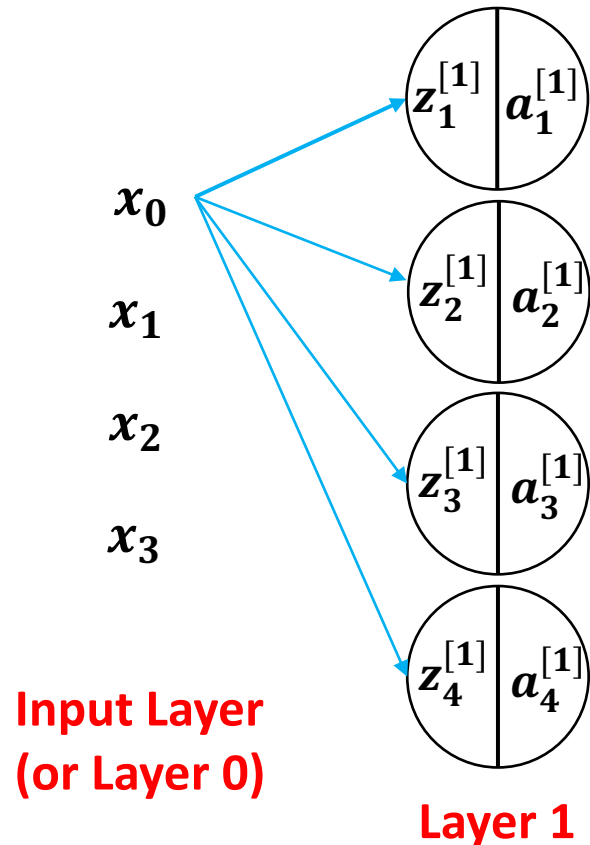
- We can construct a network of neurons (i.e., a neural network) with as many *layers*, and neurons in any layer, as needed





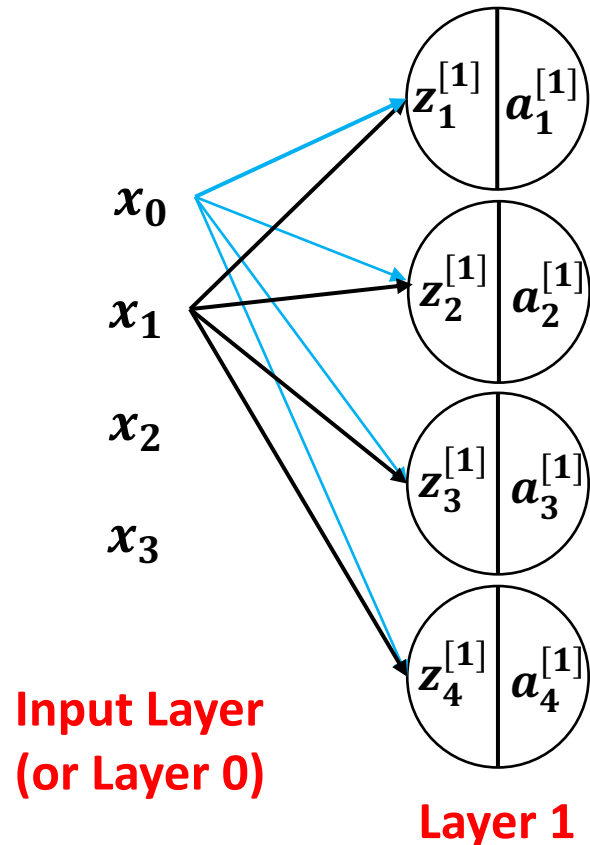
# Neural Networks

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



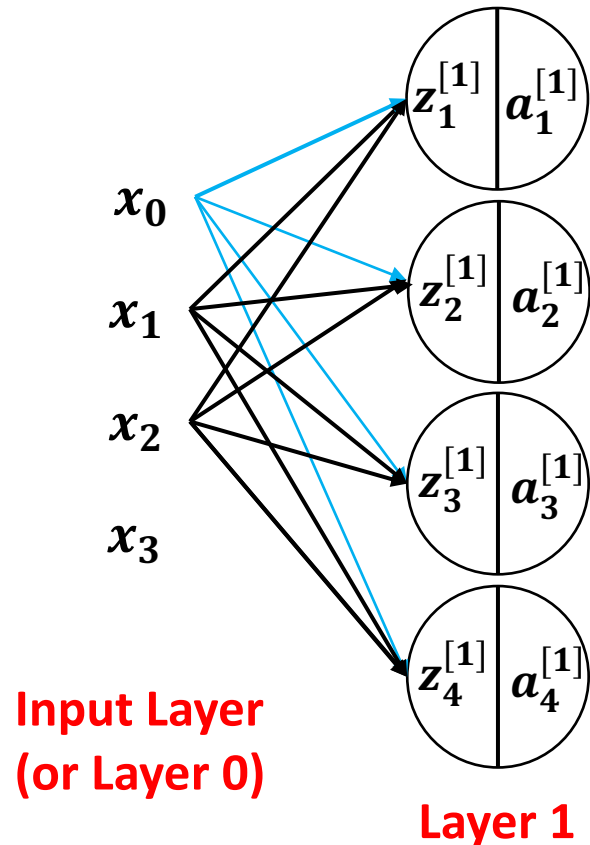
# Neural Networks

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



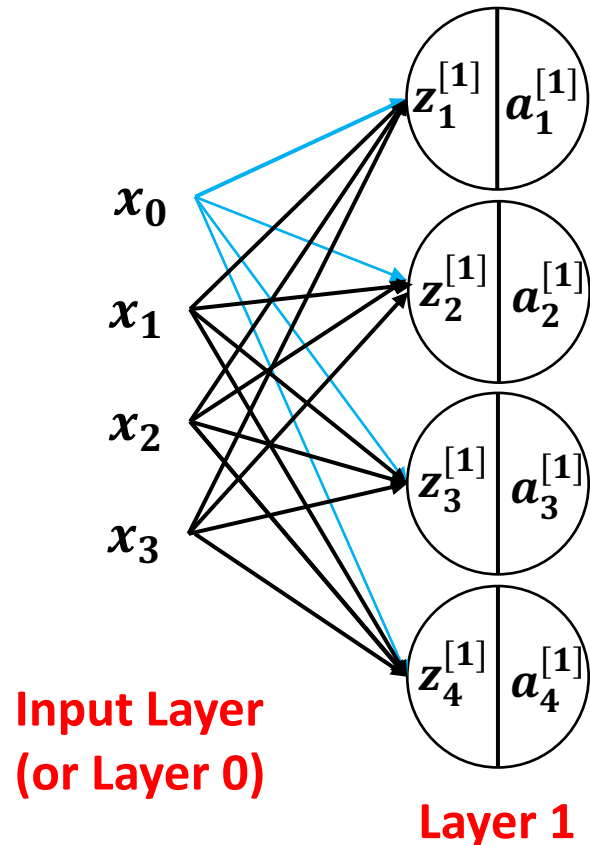
# Neural Networks

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



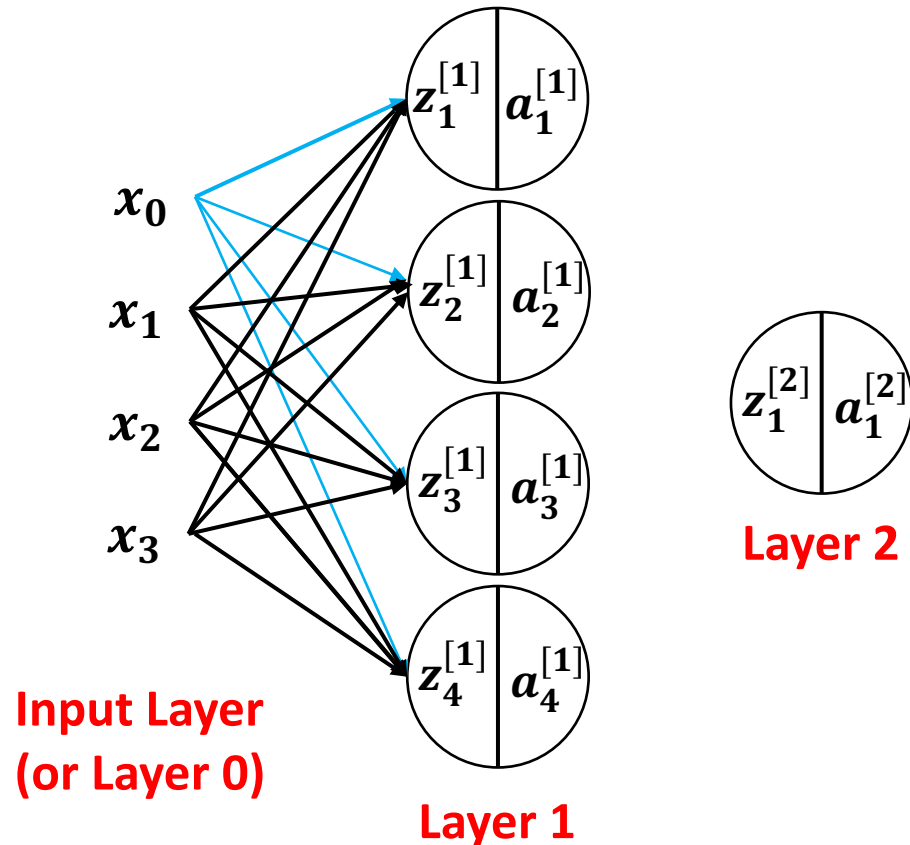
# Neural Networks

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



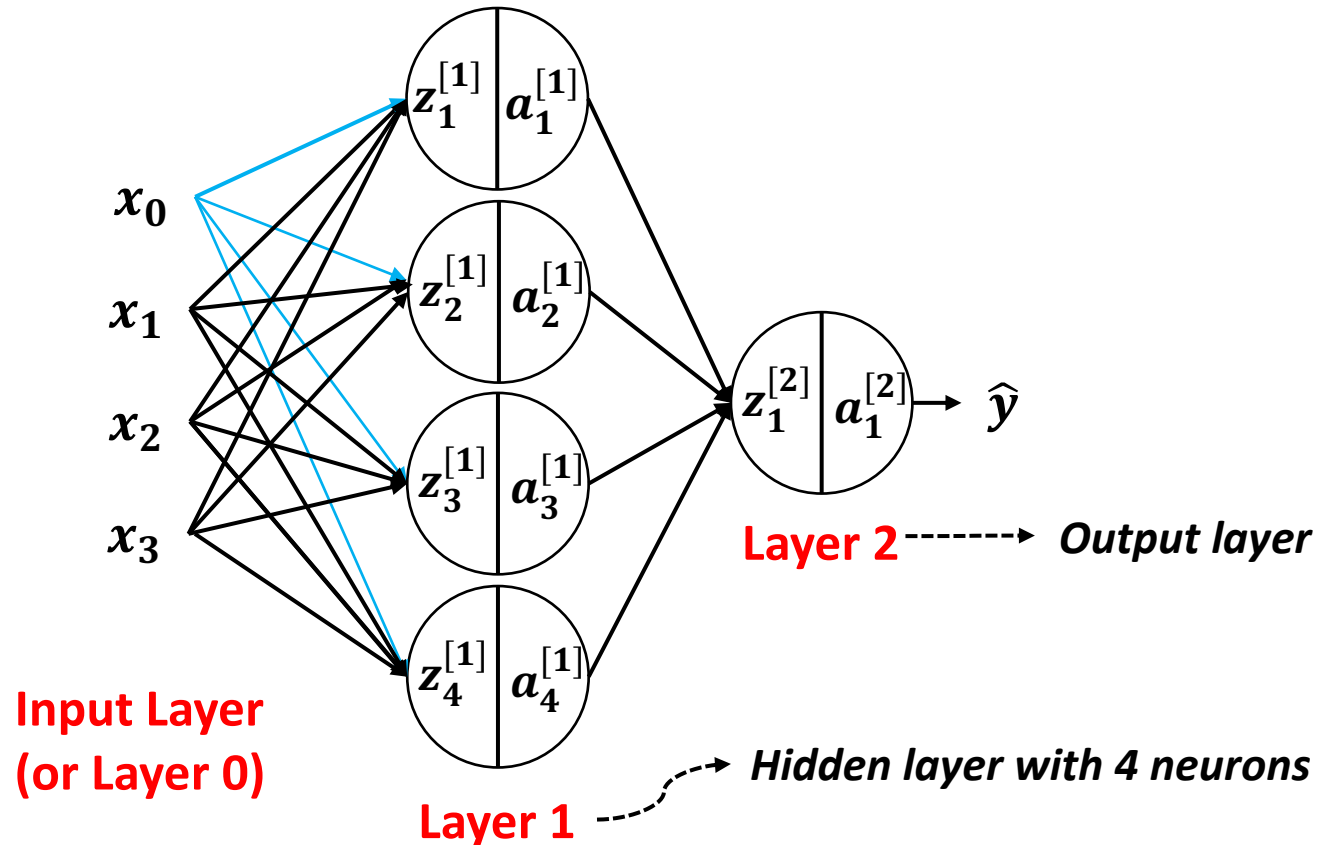
# Neural Networks

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



# Neural Networks

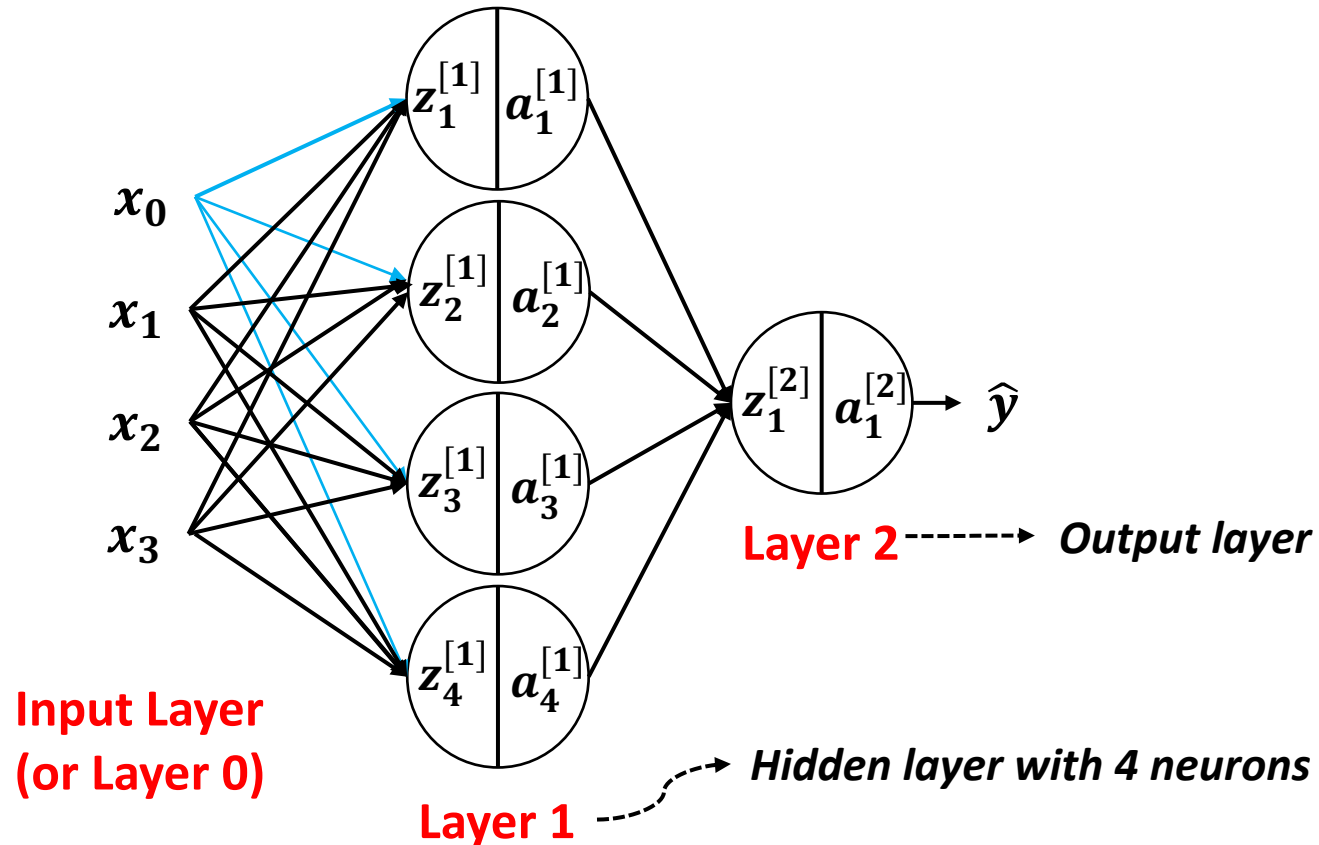
- We can construct a network of neurons (i.e., a neural network) with as many *layers*, and neurons in any layer, as needed



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

# Neural Networks

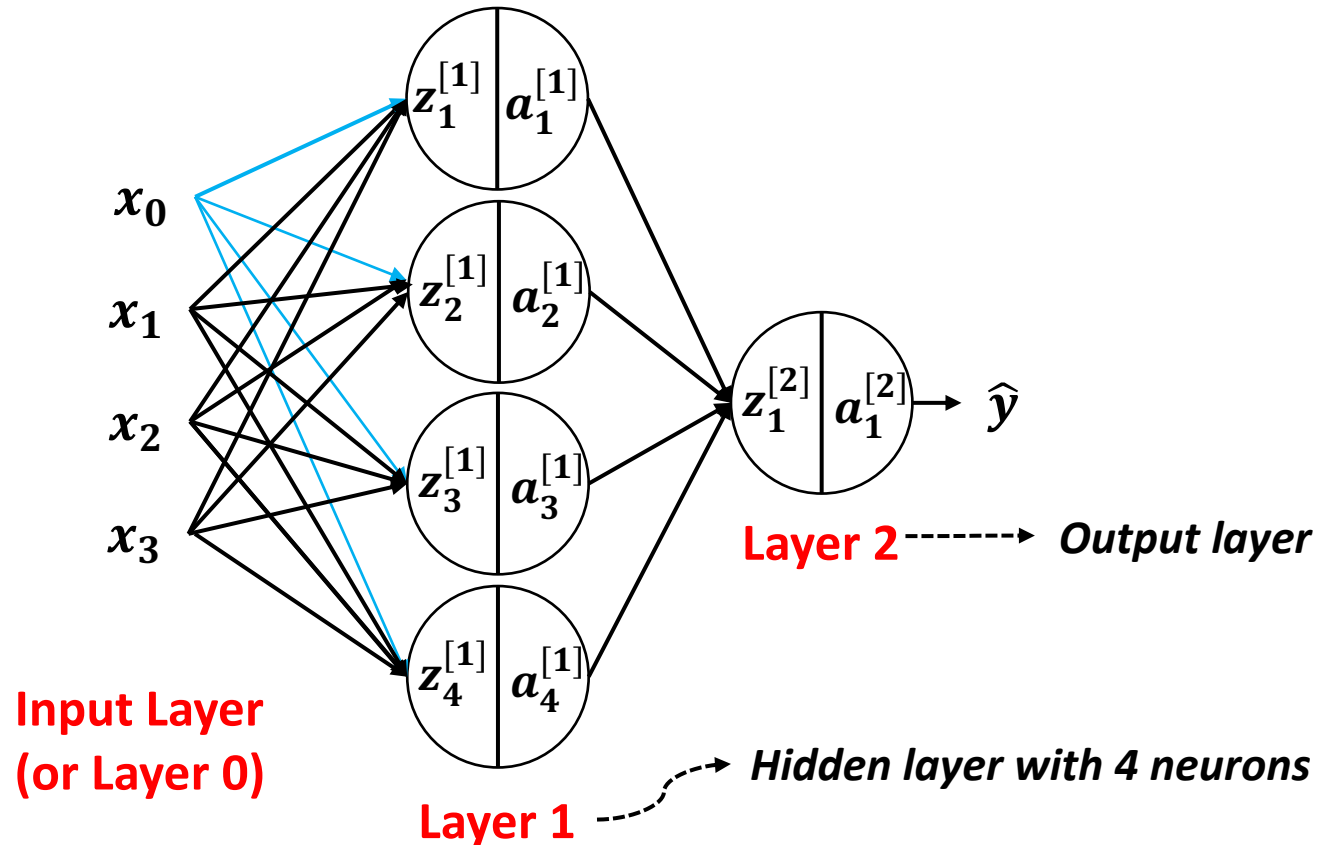
- We can construct a network of neurons (i.e., a neural network) with as many *layers*, and neurons in any layer, as needed



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

# 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

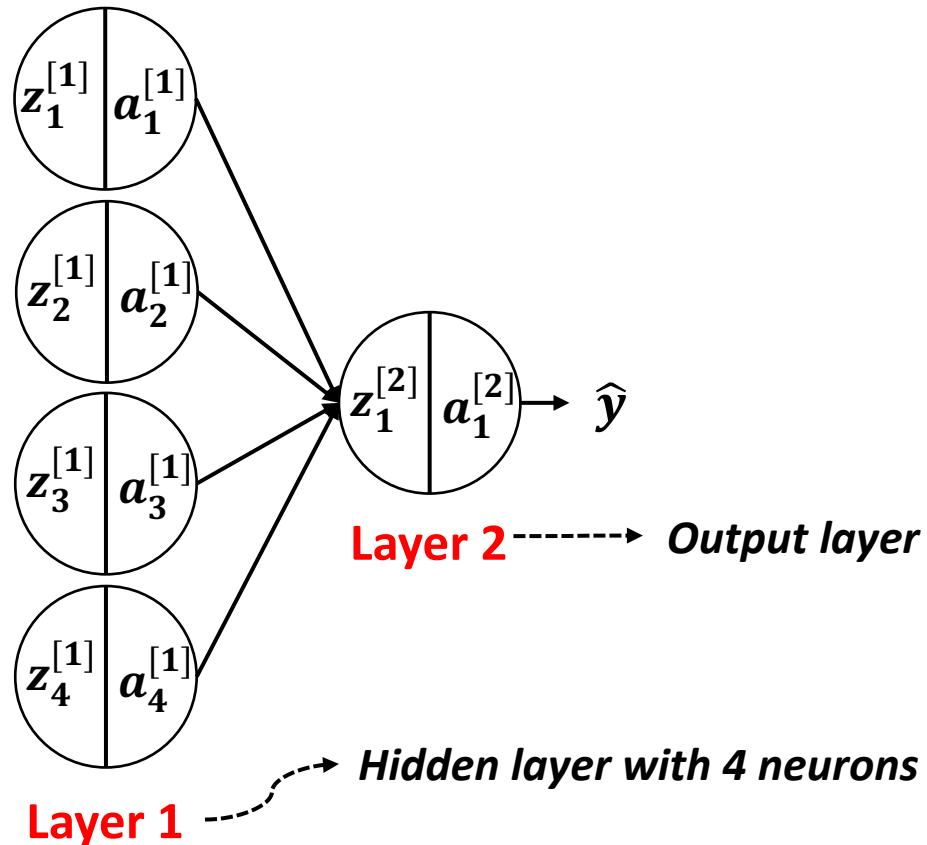


Interestingly, neural networks *learn* their own features!



# Neural Networks

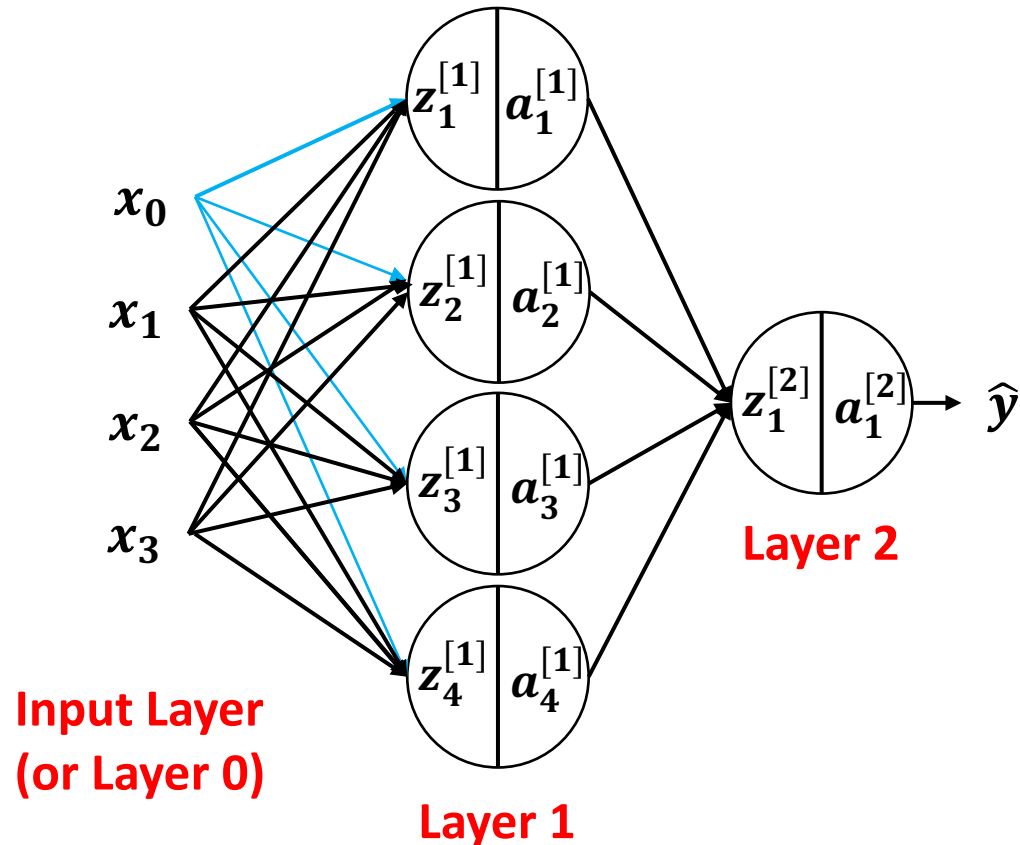
- We can construct a network of neurons (i.e., a neural network) with as many *layers*, and neurons in any layer, as needed



This looks like logistic regression, but with ***features that were learnt*** (i.e.,  $a_1^{[1]}$ ,  $a_2^{[1]}$ ,  $a_3^{[1]}$ ,  $a_4^{[1]}$ ) and ***NOT engineered by us*** (i.e.,  $x_1$ ,  $x_2$ , and  $x_3$ )

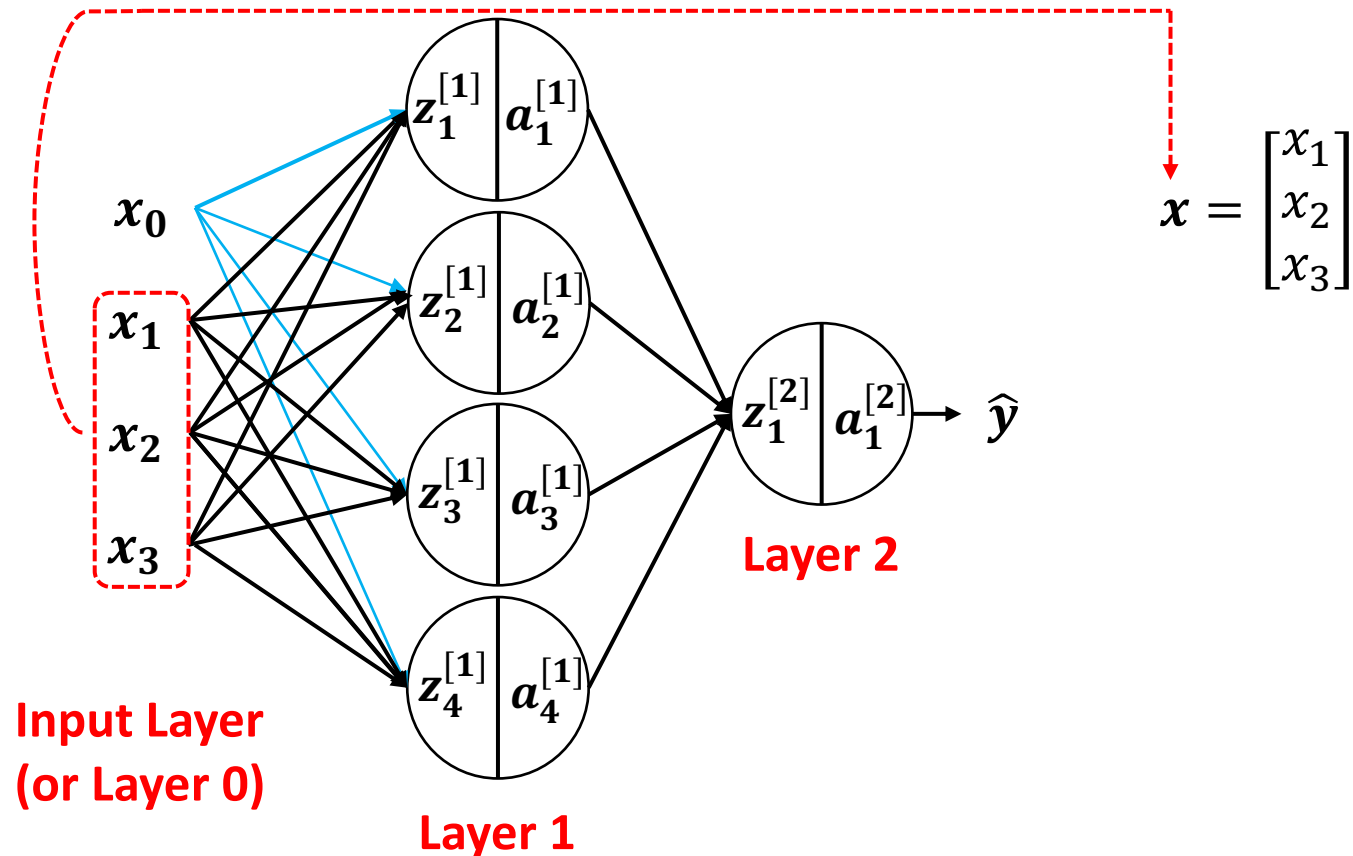
# Vectorizing Input and All Variables

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



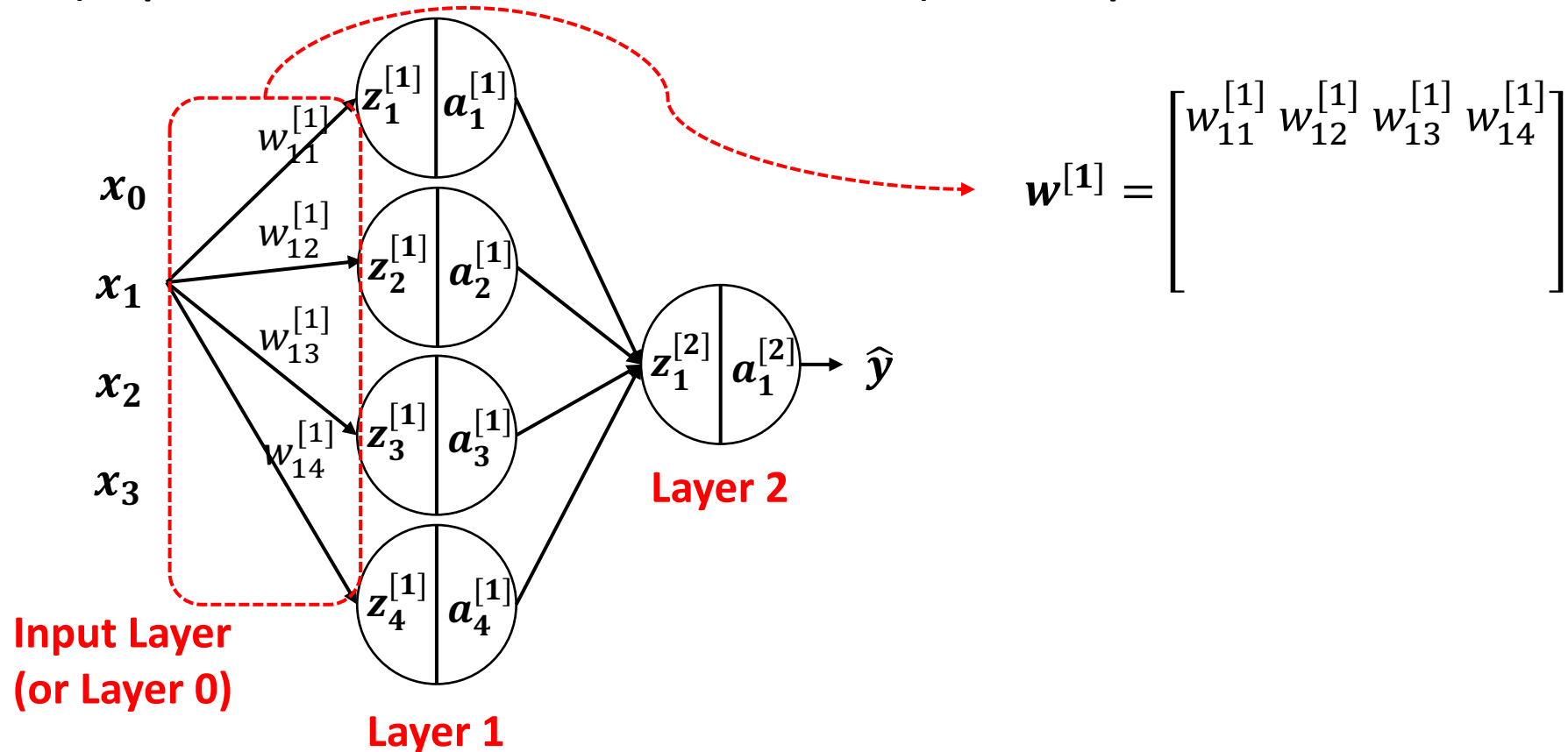
# Vectorizing Input and All Variables

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



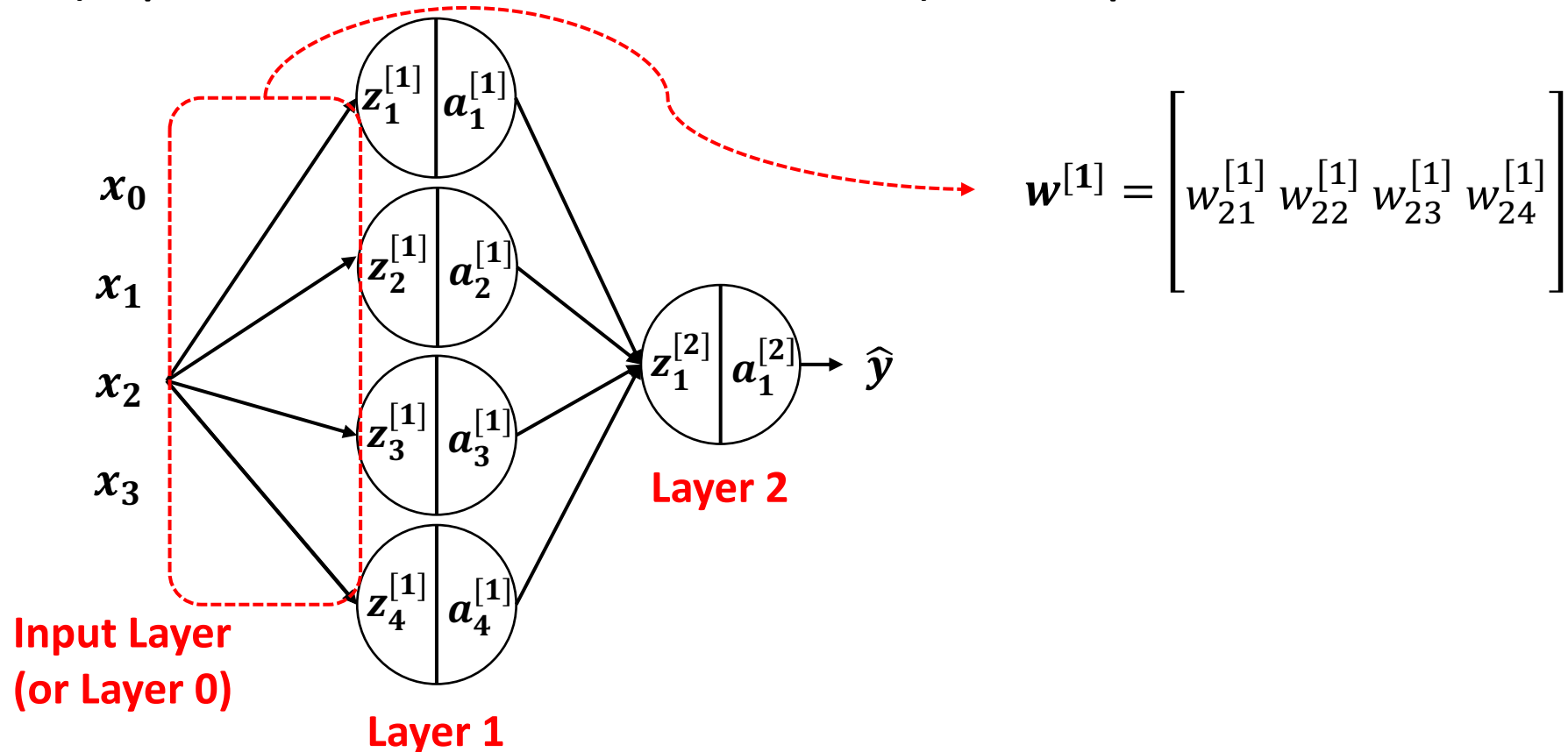
# Vectorizing Input and All Variables

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



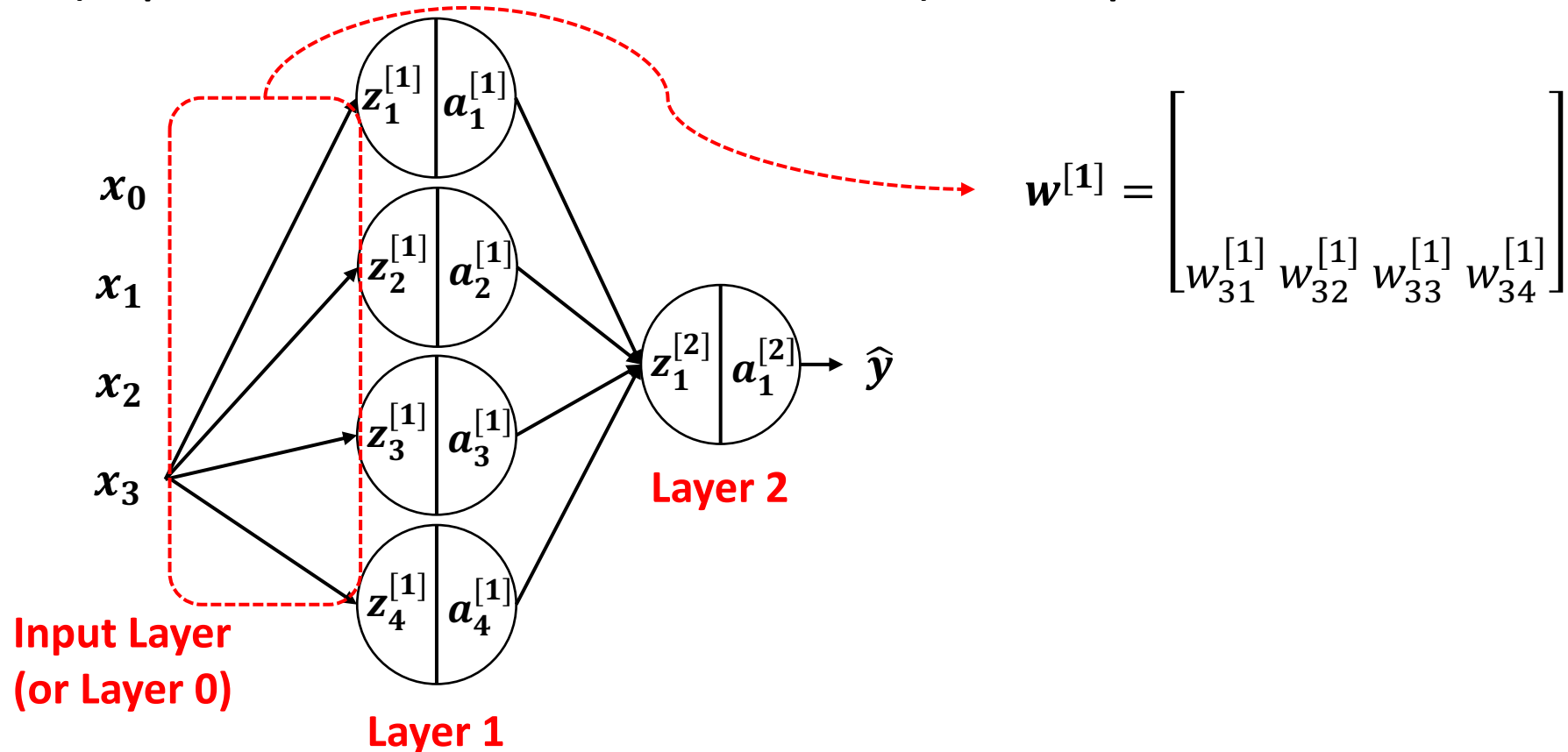
# Vectorizing Input and All Variables

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



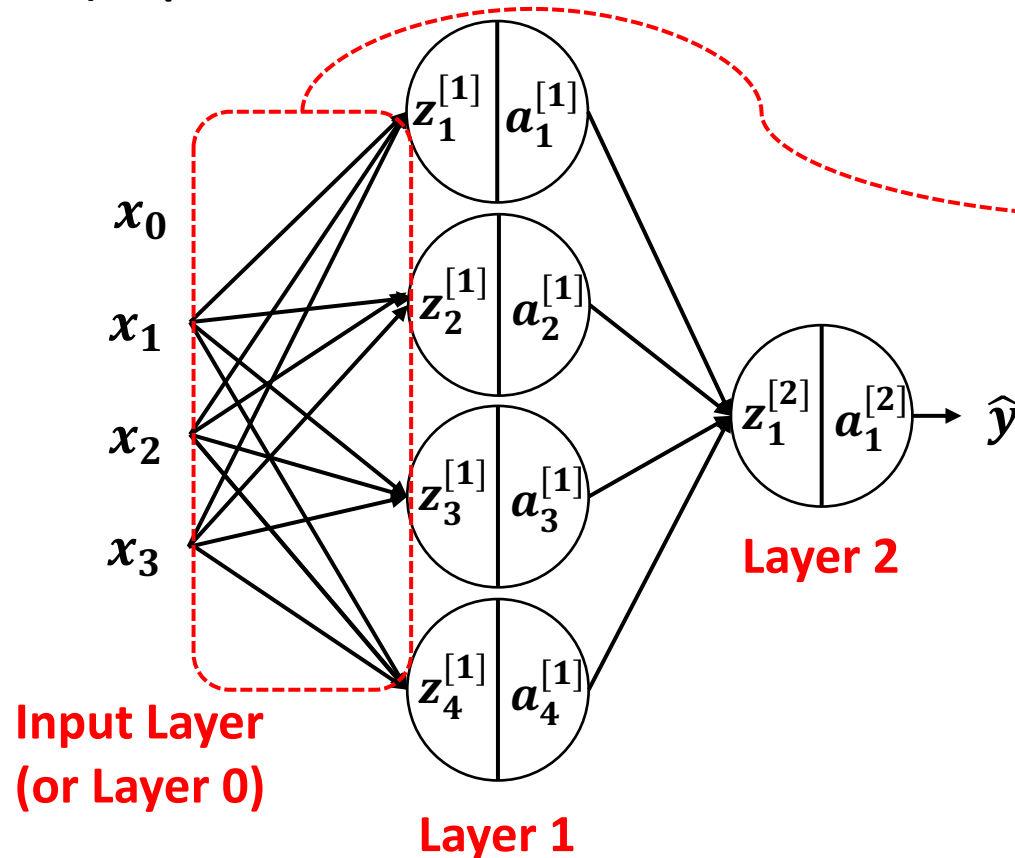
# Vectorizing Input and All Variables

- 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

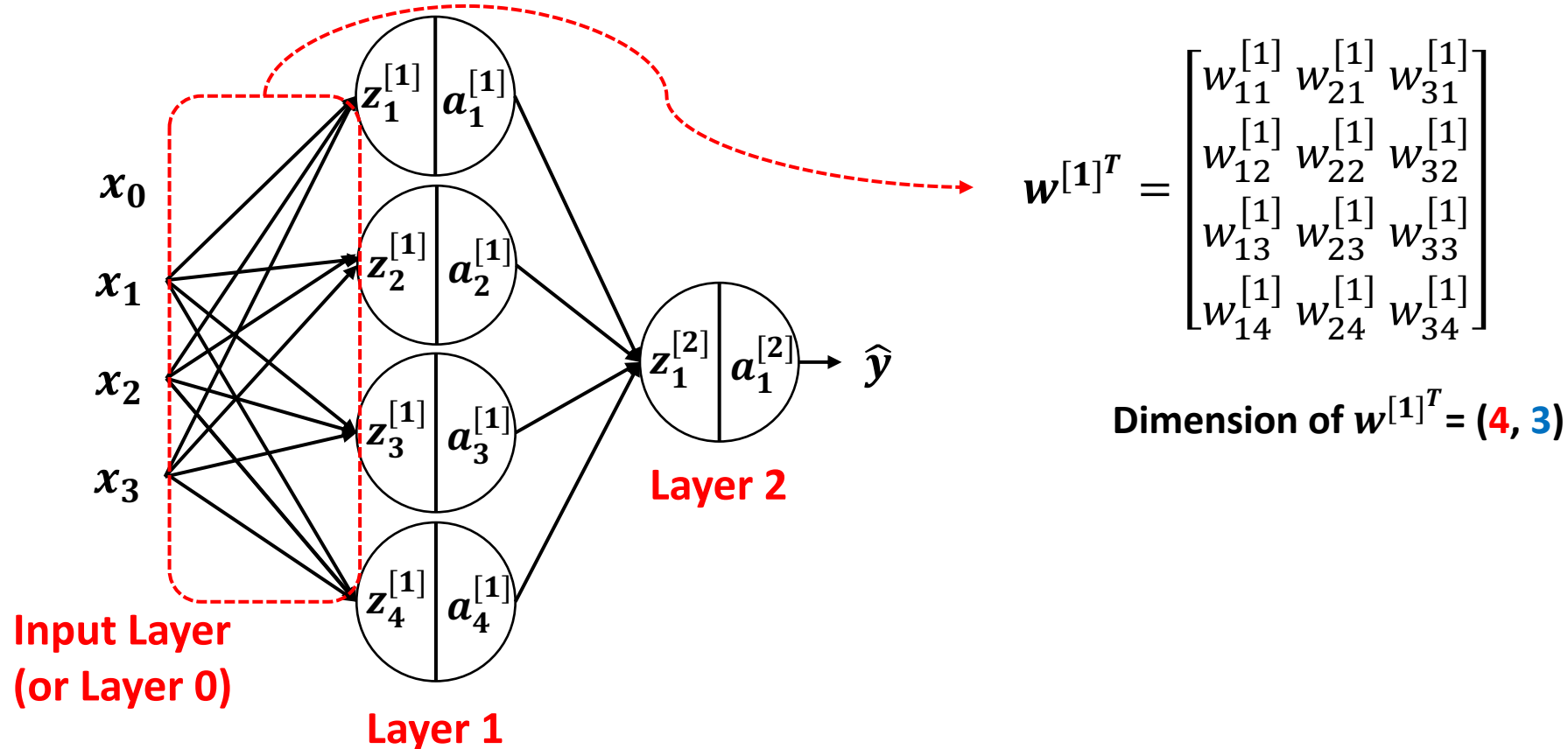


$$w^{[1]} = \begin{bmatrix} w_{11}^{[1]} & w_{12}^{[1]} & w_{13}^{[1]} & w_{14}^{[1]} \\ w_{21}^{[1]} & w_{22}^{[1]} & w_{23}^{[1]} & w_{24}^{[1]} \\ w_{31}^{[1]} & w_{32}^{[1]} & w_{33}^{[1]} & w_{34}^{[1]} \end{bmatrix}$$

Dimension of  $w^{[1]} = (3, 4)$

# Vectorizing Input and All Variables

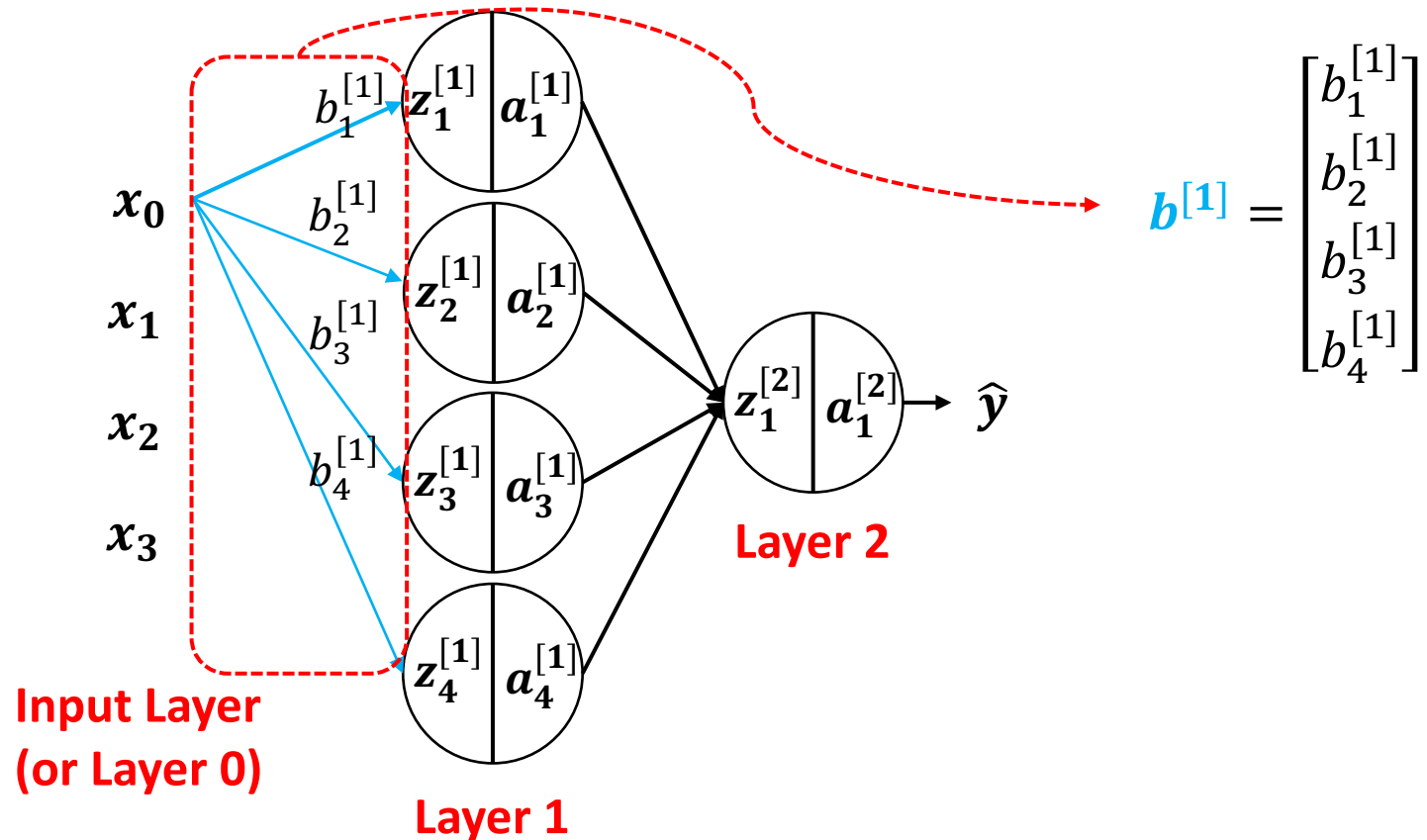
- To help develop an efficient learning algorithm, let us **vectorize** (represent in vectors & matrices) the input and the variables involved





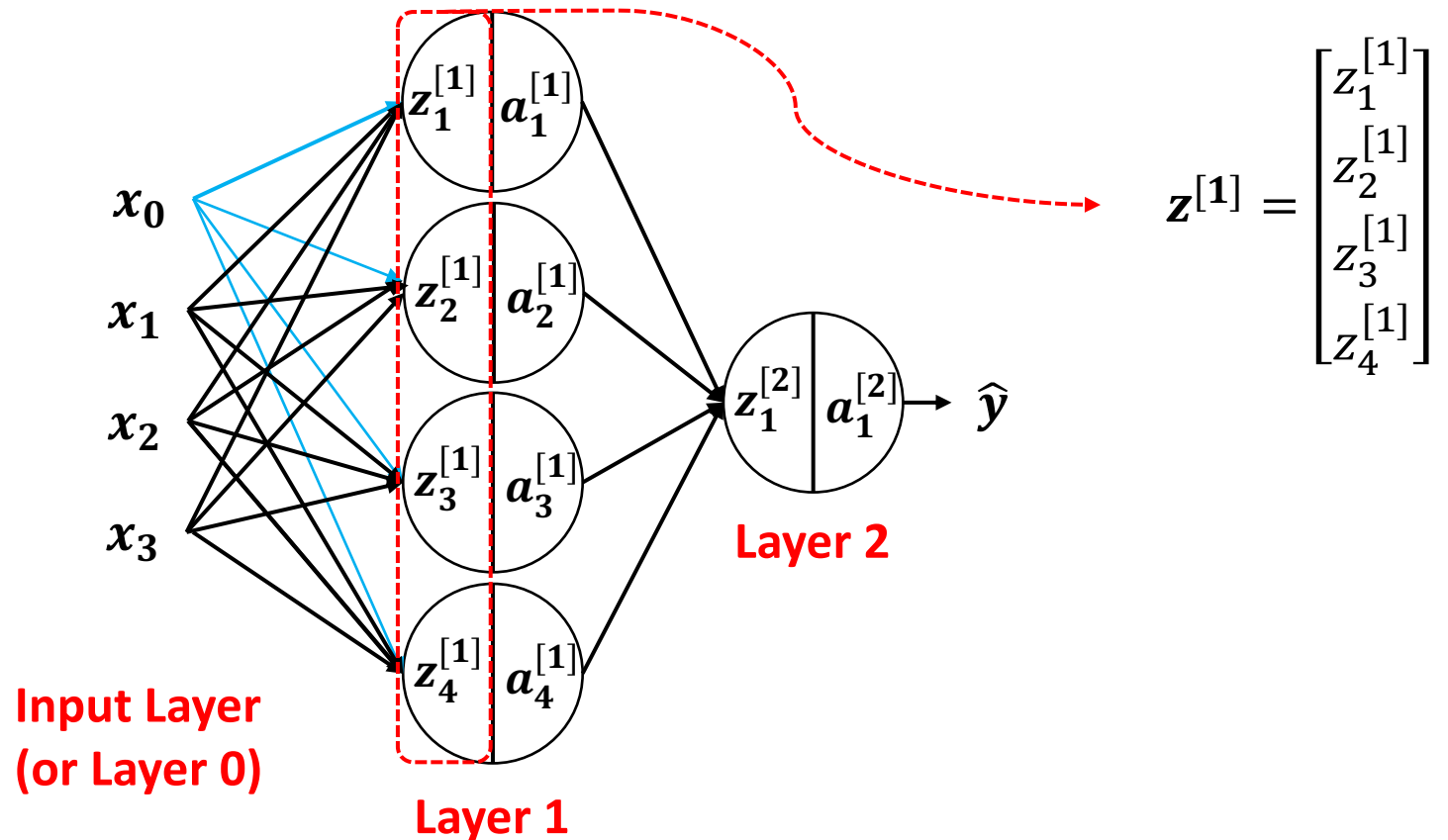
# Vectorizing Input and All Variables

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



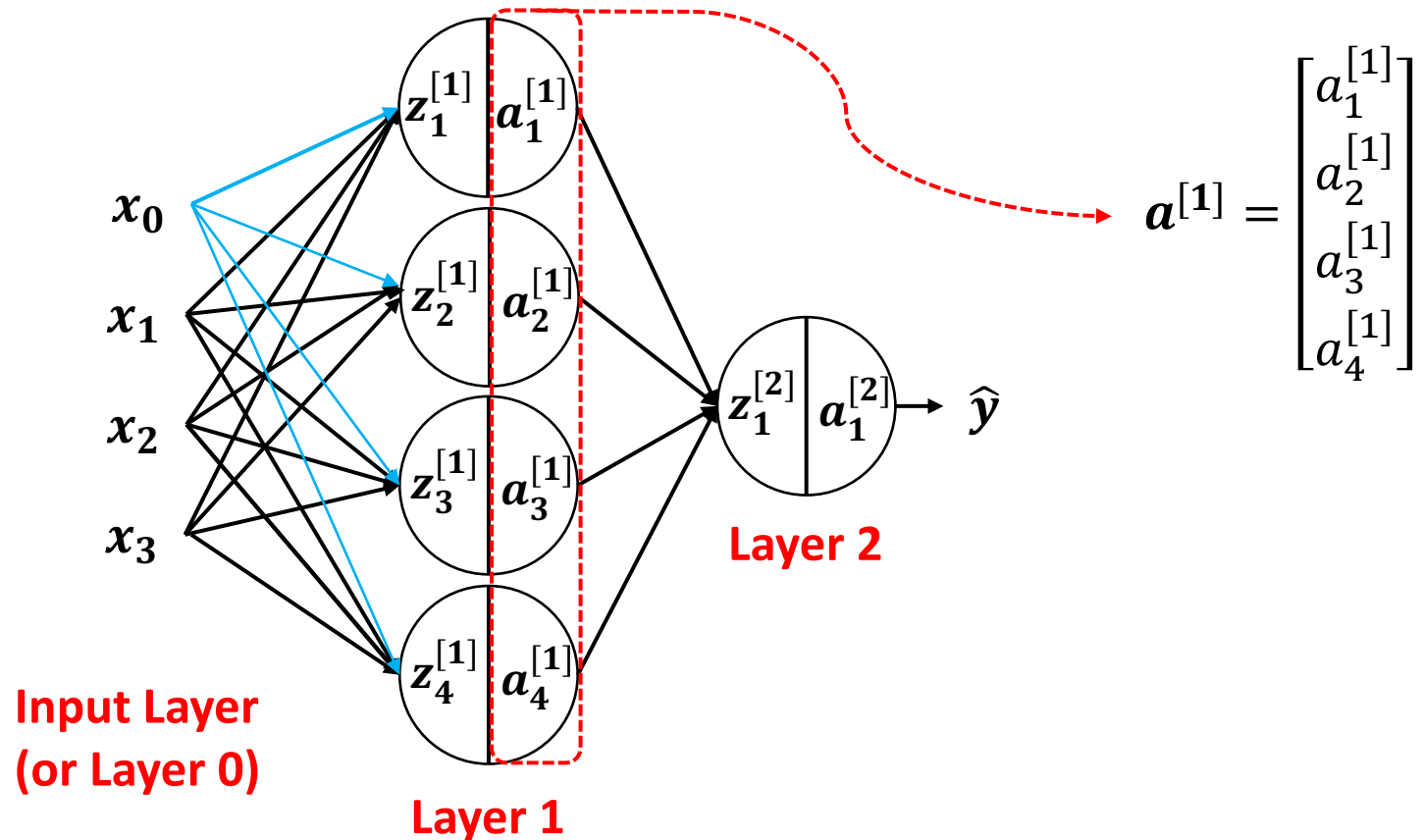
# Vectorizing Input and All Variables

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



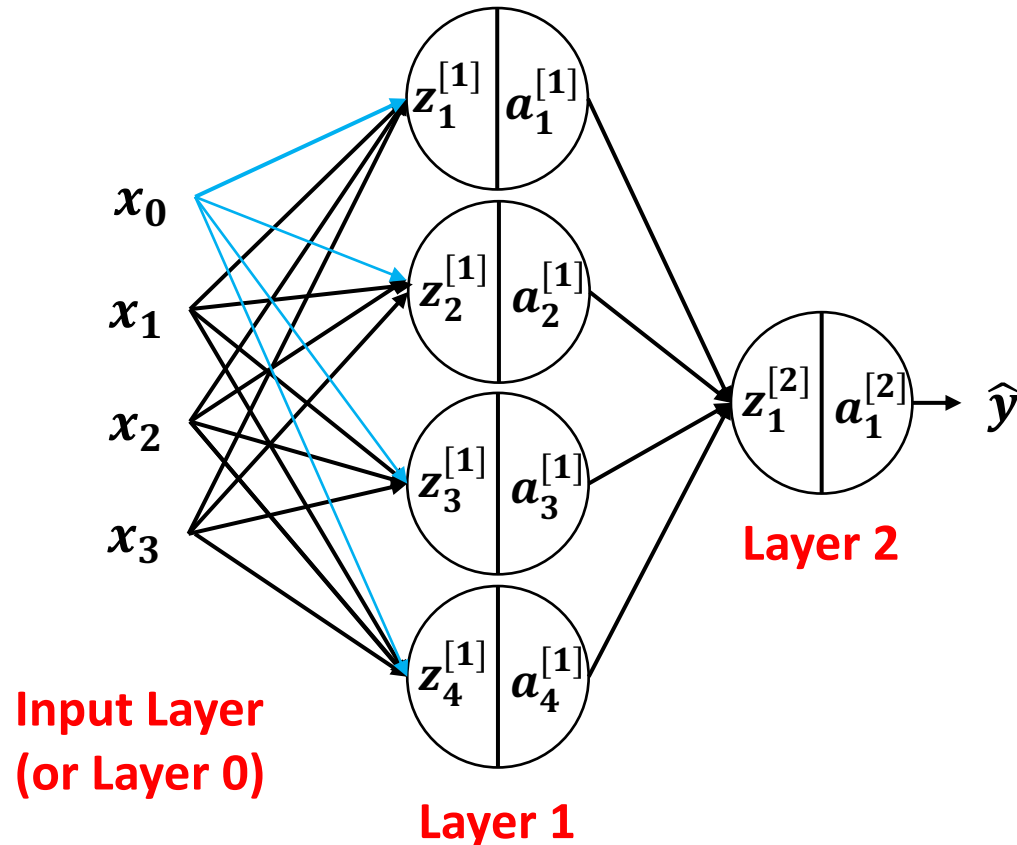
# Vectorizing Input and All Variables

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



# 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

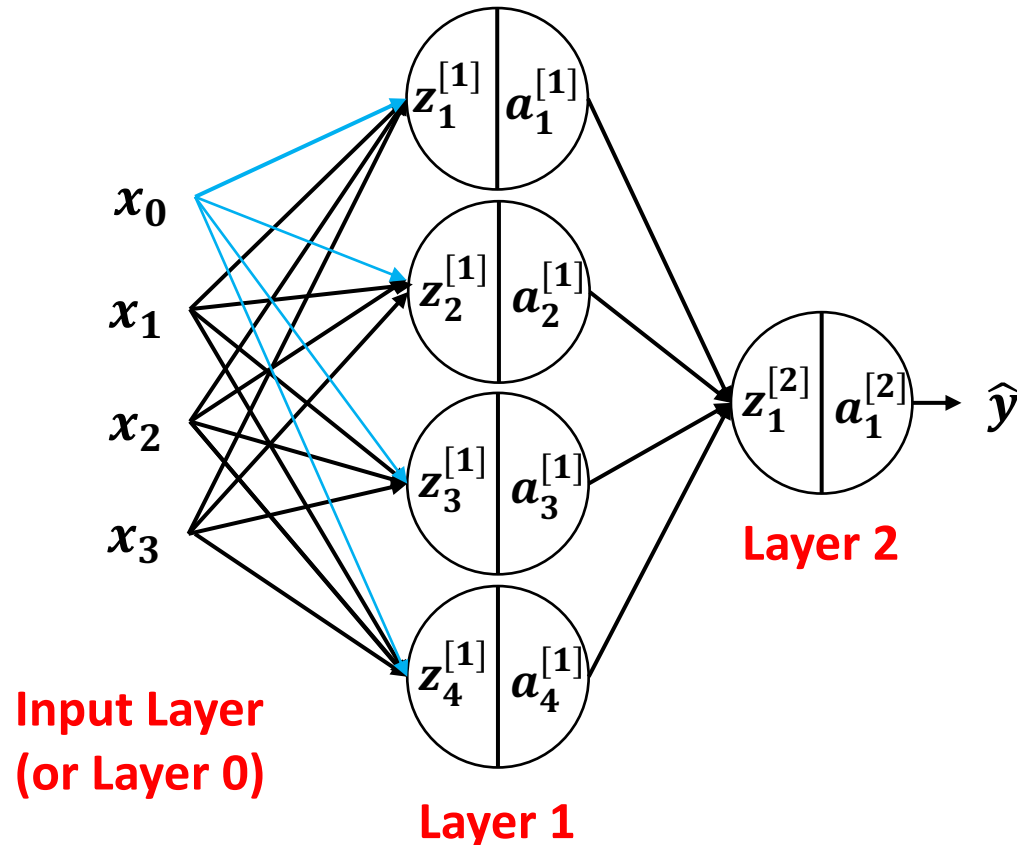


$$\mathbf{z}^{[1]} = \mathbf{w}^{[1]T} \mathbf{x} + \mathbf{b}^{[1]}$$

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

# 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

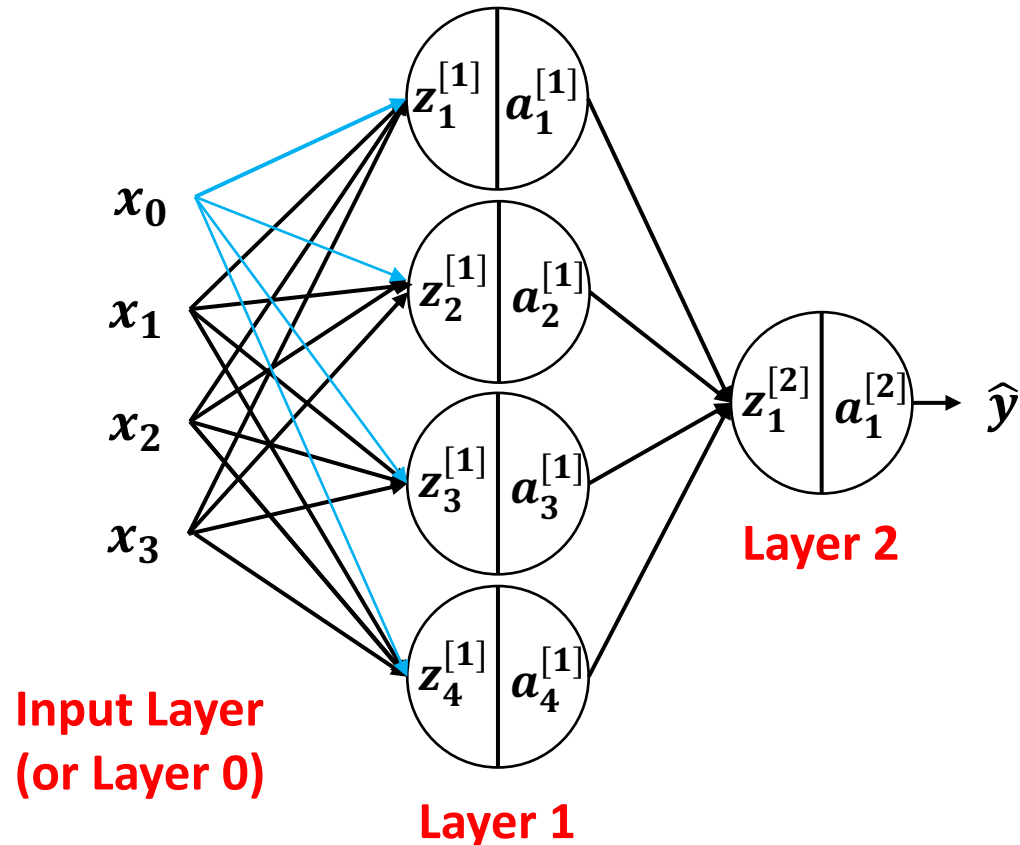


$$\mathbf{z}^{[1]} = \mathbf{w}^{[1]T} \mathbf{x} + \mathbf{b}^{[1]}$$

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

# 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



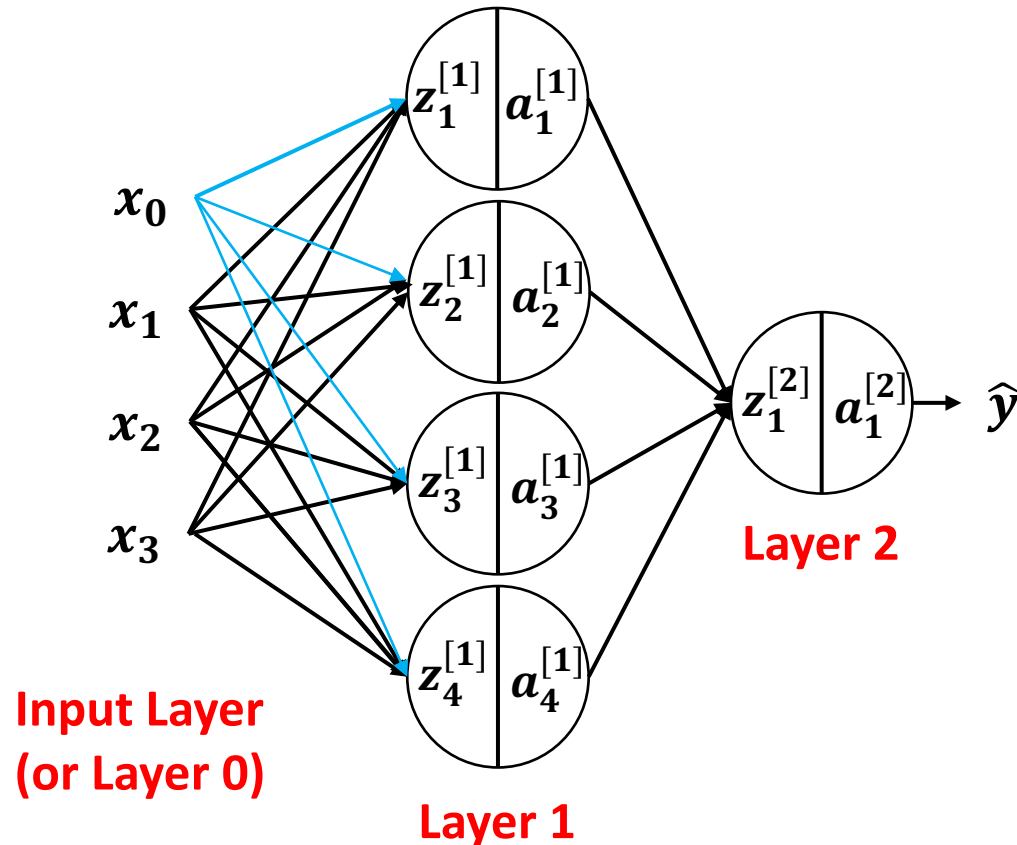
$$\mathbf{z}^{[1]} = \mathbf{w}^{[1]T} \mathbf{x} + \mathbf{b}^{[1]}$$

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

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

# 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



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

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

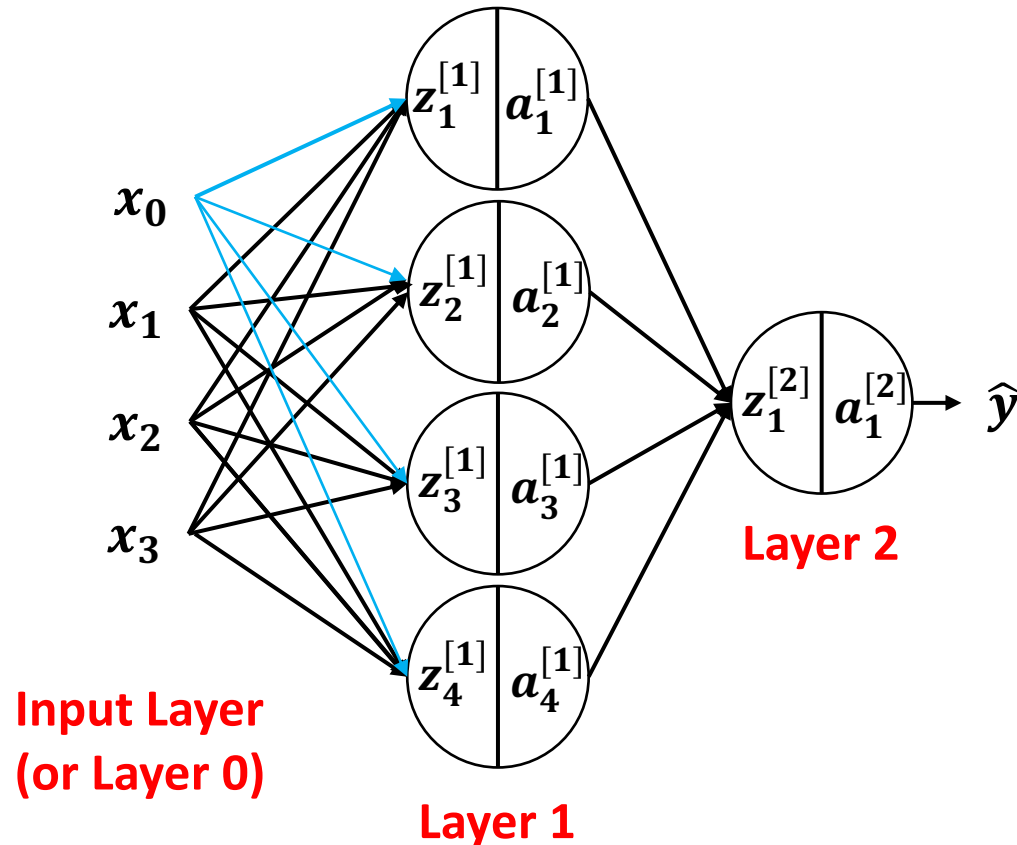
A green dashed arrow points from the first row of the vector equation to the  $z_1^{[1]}$  label below.

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

$z_1^{[1]}$

# 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



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

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

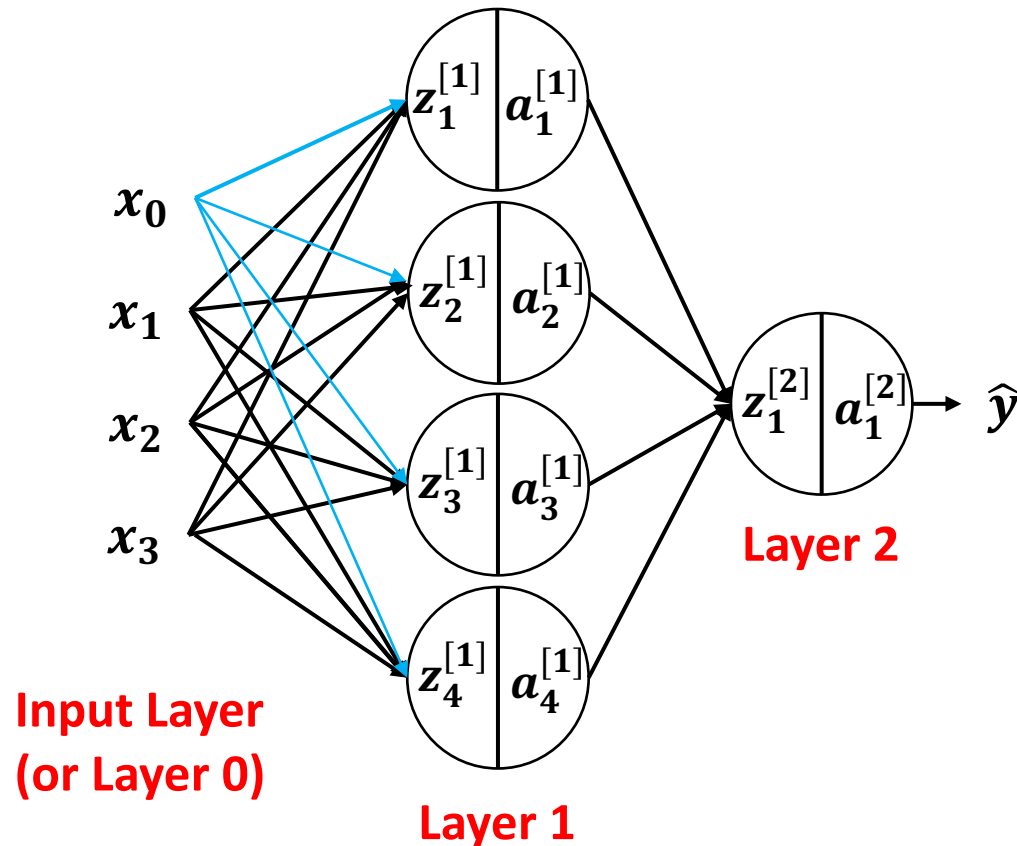
$$\mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]})$$

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



# 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



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

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

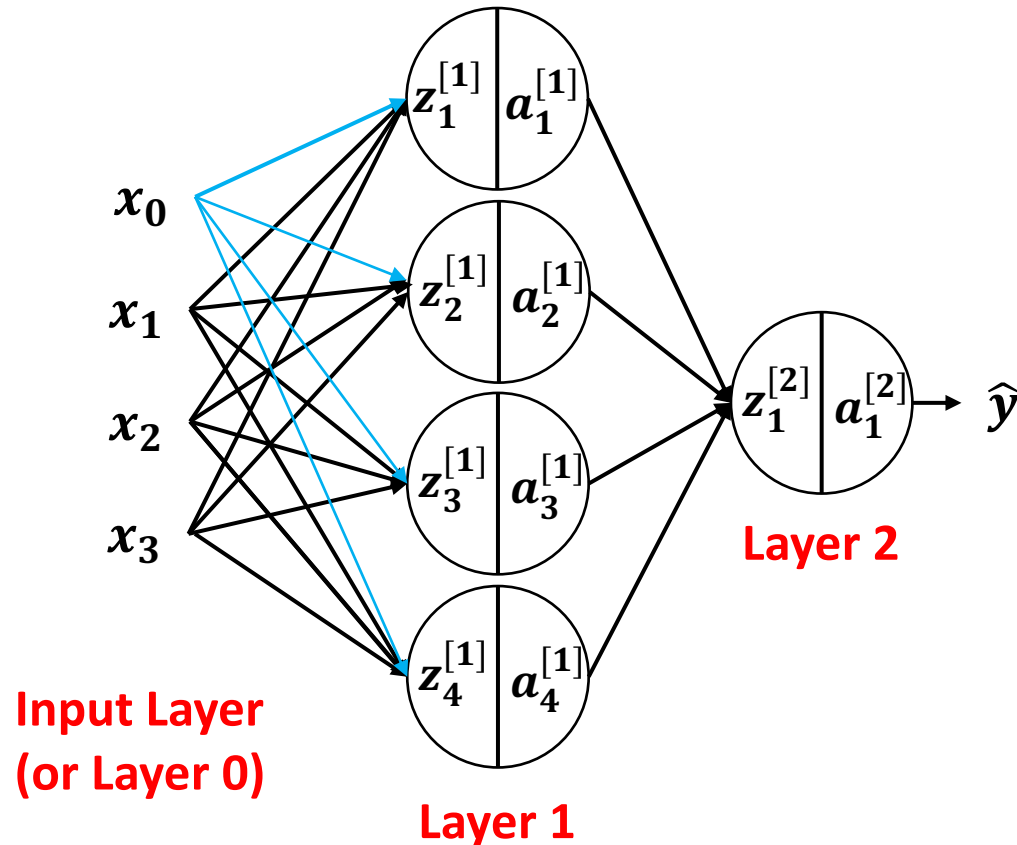
A green dashed arrow points from the second row of the vector  $\mathbf{z}^{[1]}$  to the label  $z_2^{[1]}$  at the bottom right.

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

$z_2^{[1]}$

# 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



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

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

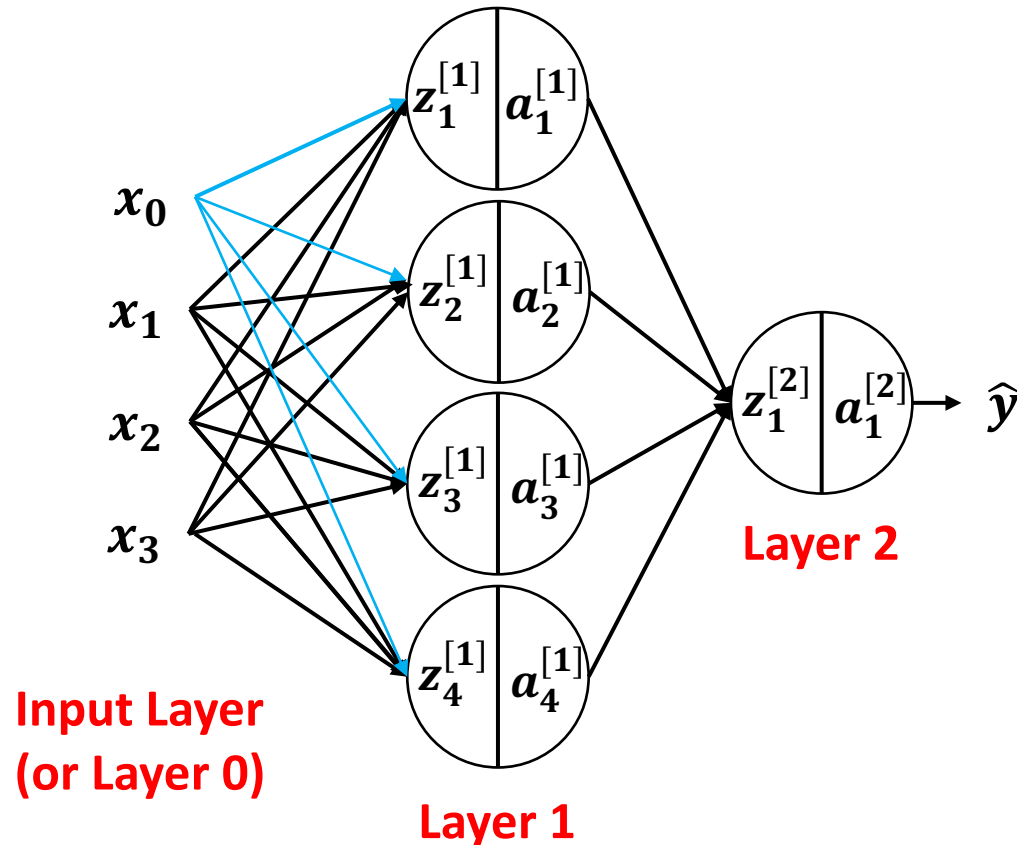
A green dashed arrow points from the second row of the matrix to the equation  $a_2^{[1]} = \sigma(z_2^{[1]})$  below.

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

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

# 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



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

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

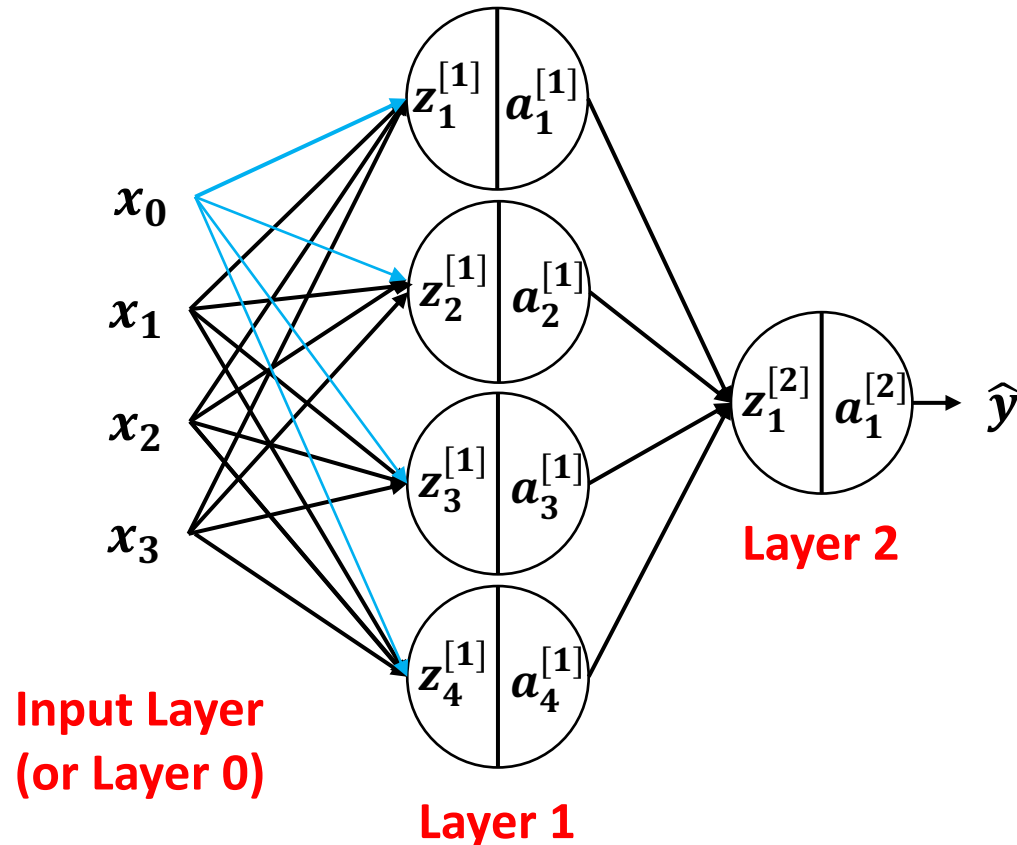
A green box highlights the third row of the vector, which corresponds to  $z_3^{[1]}$  as indicated by a dashed green arrow.

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

$z_3^{[1]}$

# 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



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

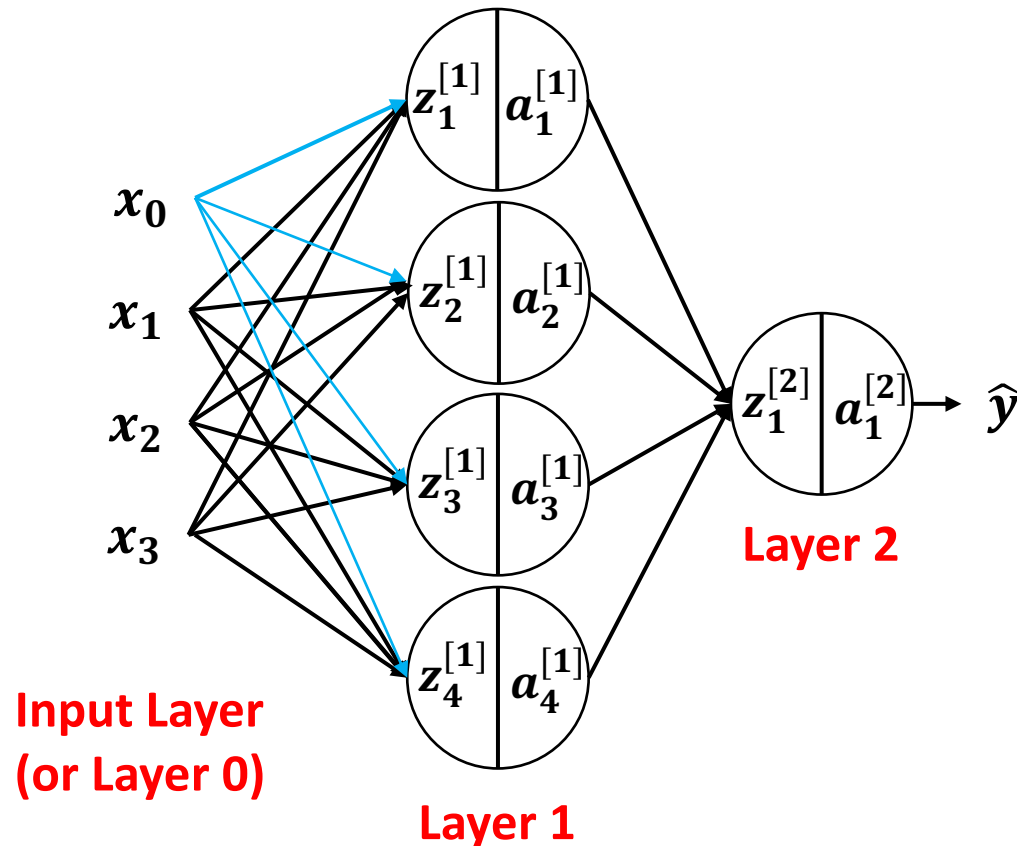
$$= \begin{bmatrix} w_{11}^{[1]} x_1 + w_{21}^{[1]} x_2 + w_{31}^{[1]} x_3 + b_1^{[1]} \\ w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2 + w_{32}^{[1]} x_3 + b_2^{[1]} \\ w_{13}^{[1]} x_1 + w_{23}^{[1]} x_2 + w_{33}^{[1]} x_3 + b_3^{[1]} \\ w_{14}^{[1]} x_1 + w_{24}^{[1]} x_2 + w_{34}^{[1]} x_3 + b_4^{[1]} \end{bmatrix}$$

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

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

# 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



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

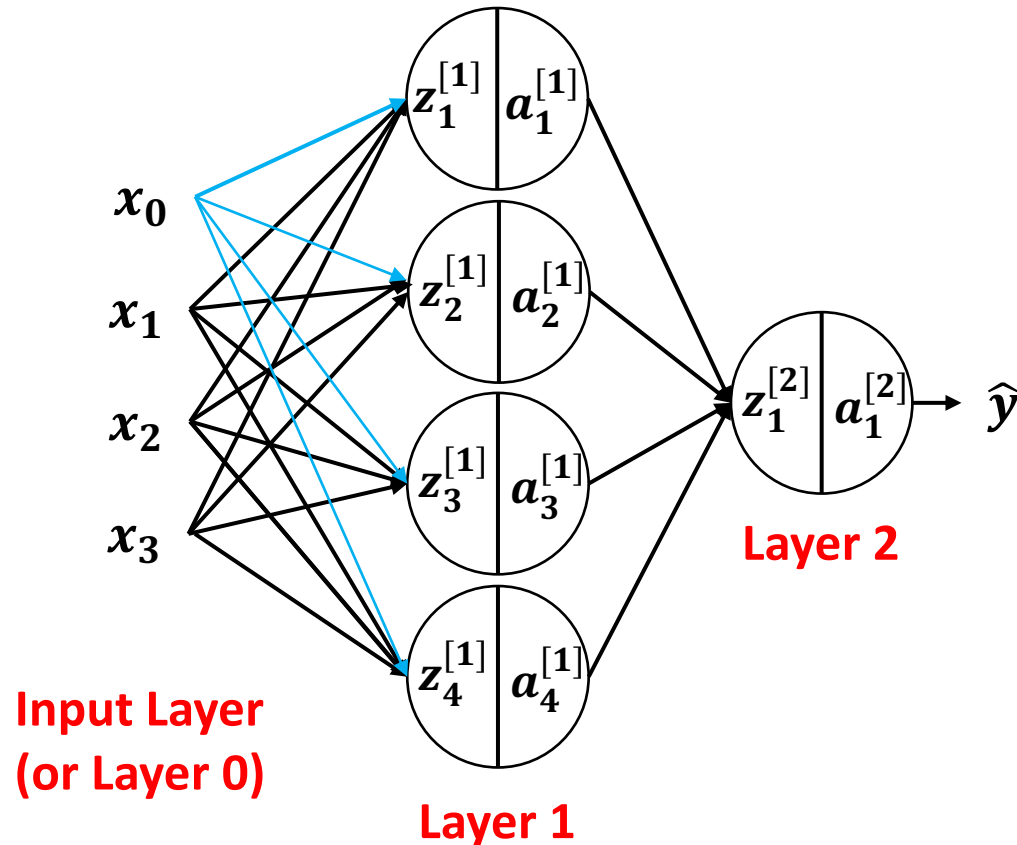
$$= \begin{bmatrix} w_{11}^{[1]} x_1 + w_{21}^{[1]} x_2 + w_{31}^{[1]} x_3 + b_1^{[1]} \\ w_{12}^{[1]} x_1 + w_{22}^{[1]} x_2 + w_{32}^{[1]} x_3 + b_2^{[1]} \\ w_{13}^{[1]} x_1 + w_{23}^{[1]} x_2 + w_{33}^{[1]} x_3 + b_3^{[1]} \\ w_{14}^{[1]} x_1 + w_{24}^{[1]} x_2 + w_{34}^{[1]} x_3 + b_4^{[1]} \end{bmatrix}$$

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

$$\mathbf{z}_4^{[1]}$$

# 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



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

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

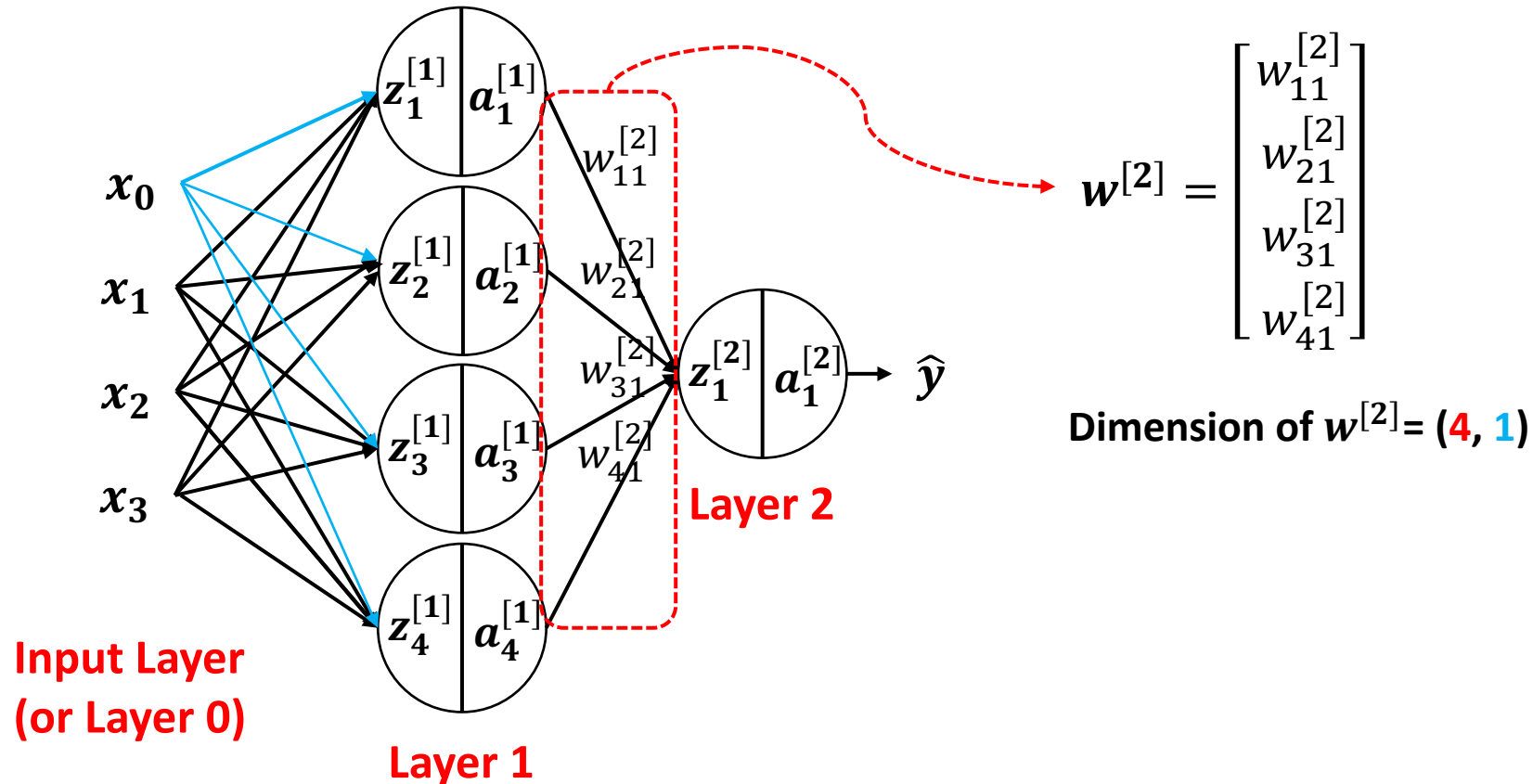
A green dashed arrow points from the bottom row of the vector equation to the activation equation below.

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

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

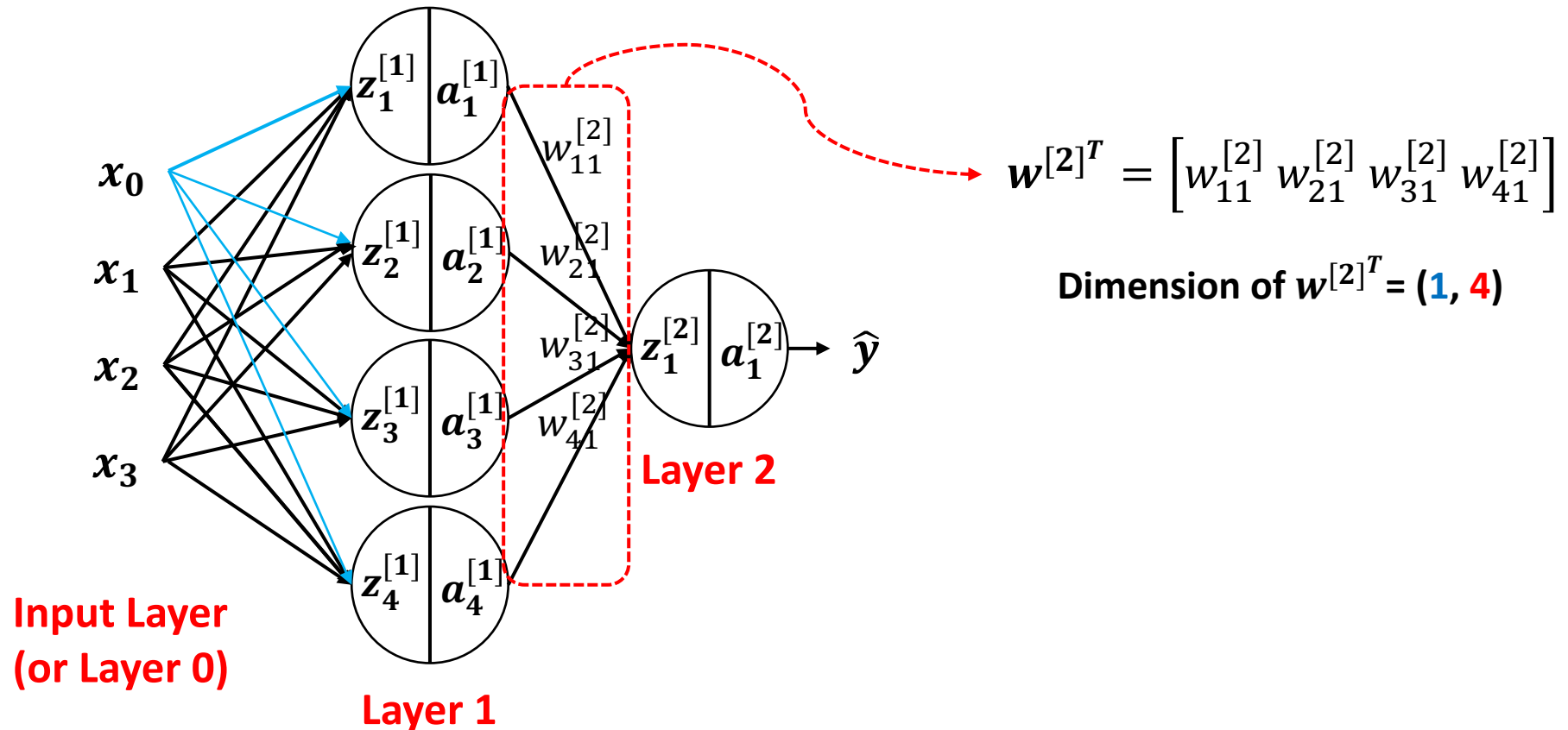
# Vectorizing Input and All Variables

- 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

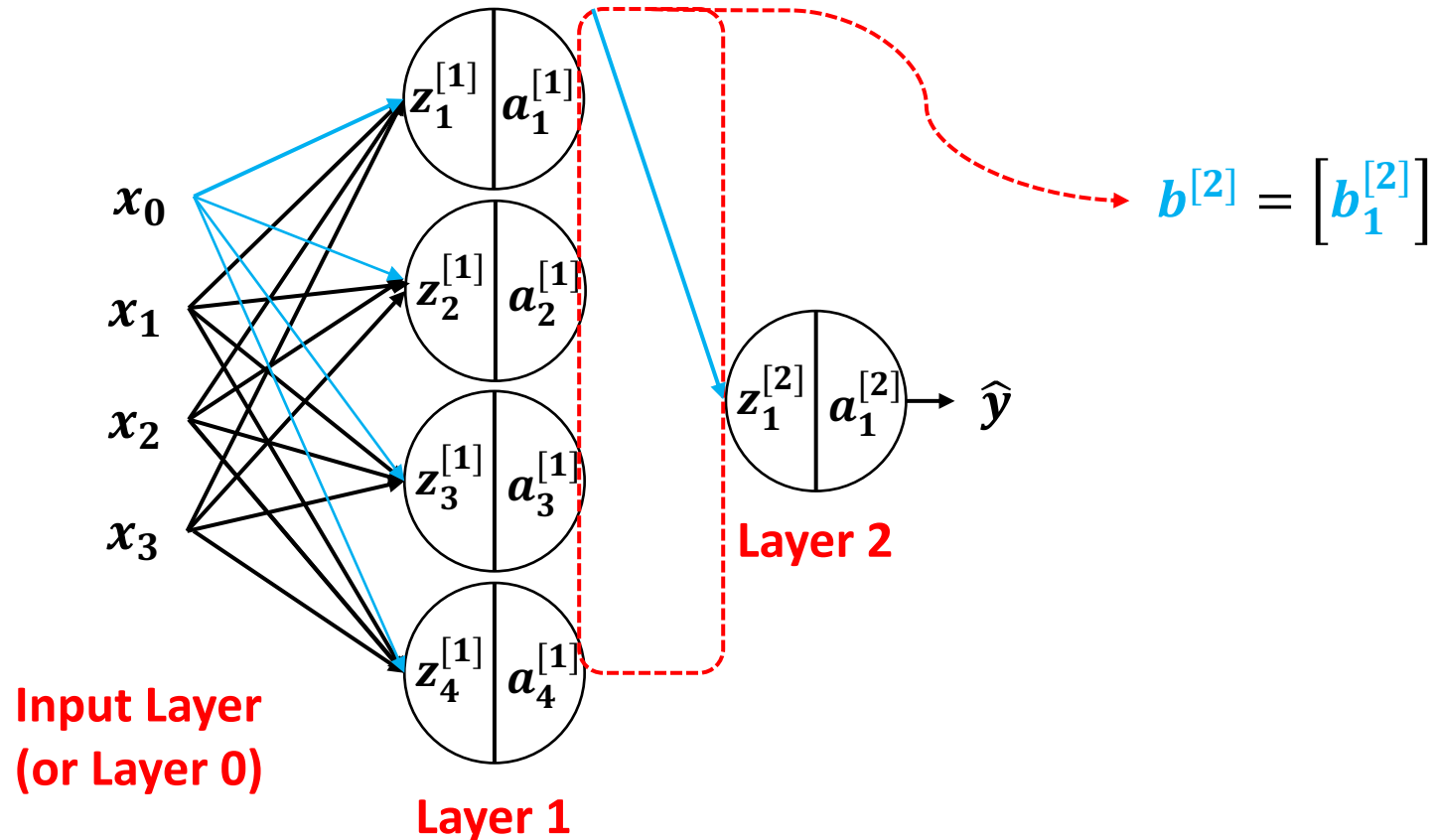
- To help develop an efficient learning algorithm, let us **vectorize** (represent in vectors & matrices) the input and the variables involved





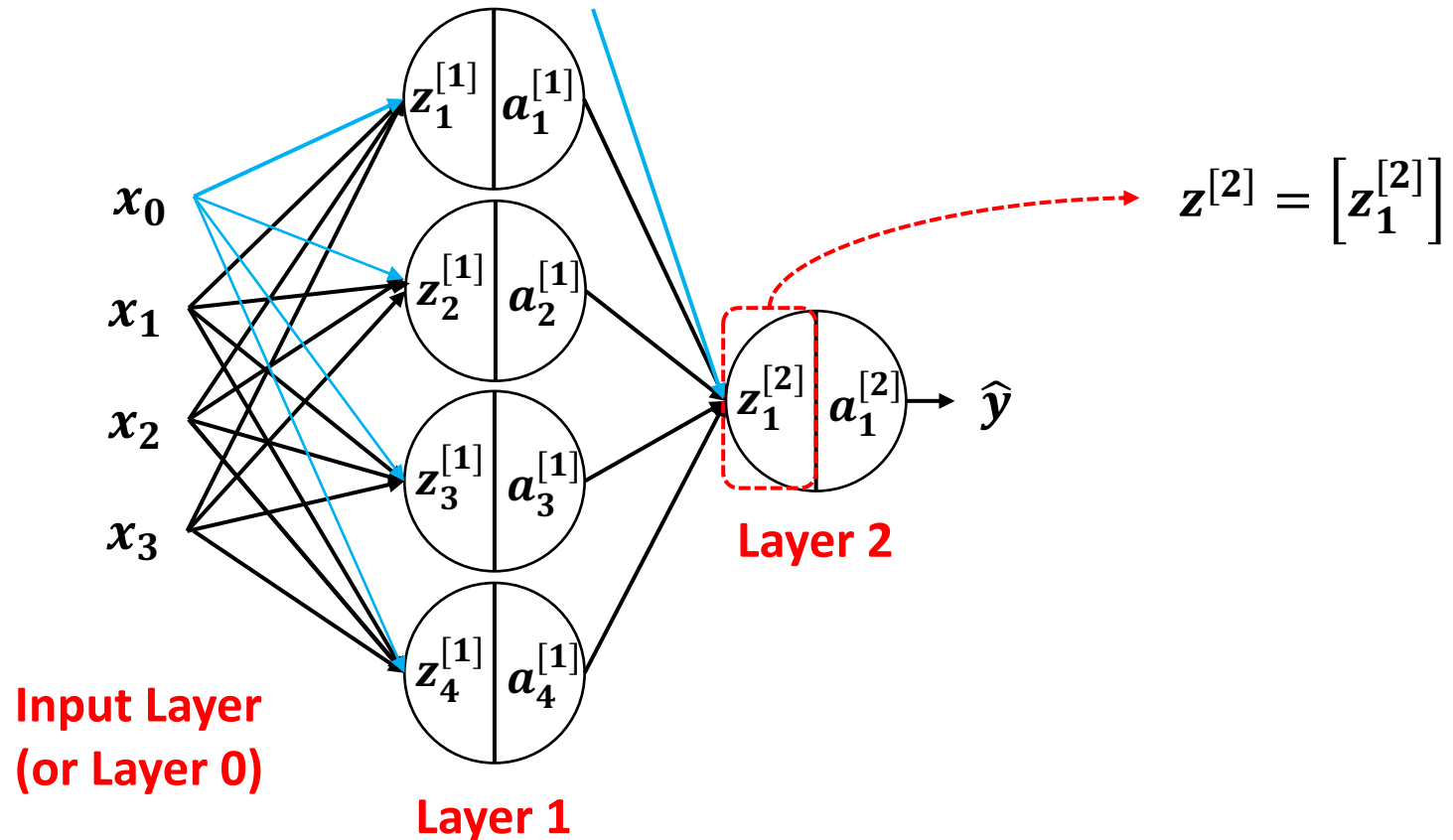
# Vectorizing Input and All Variables

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



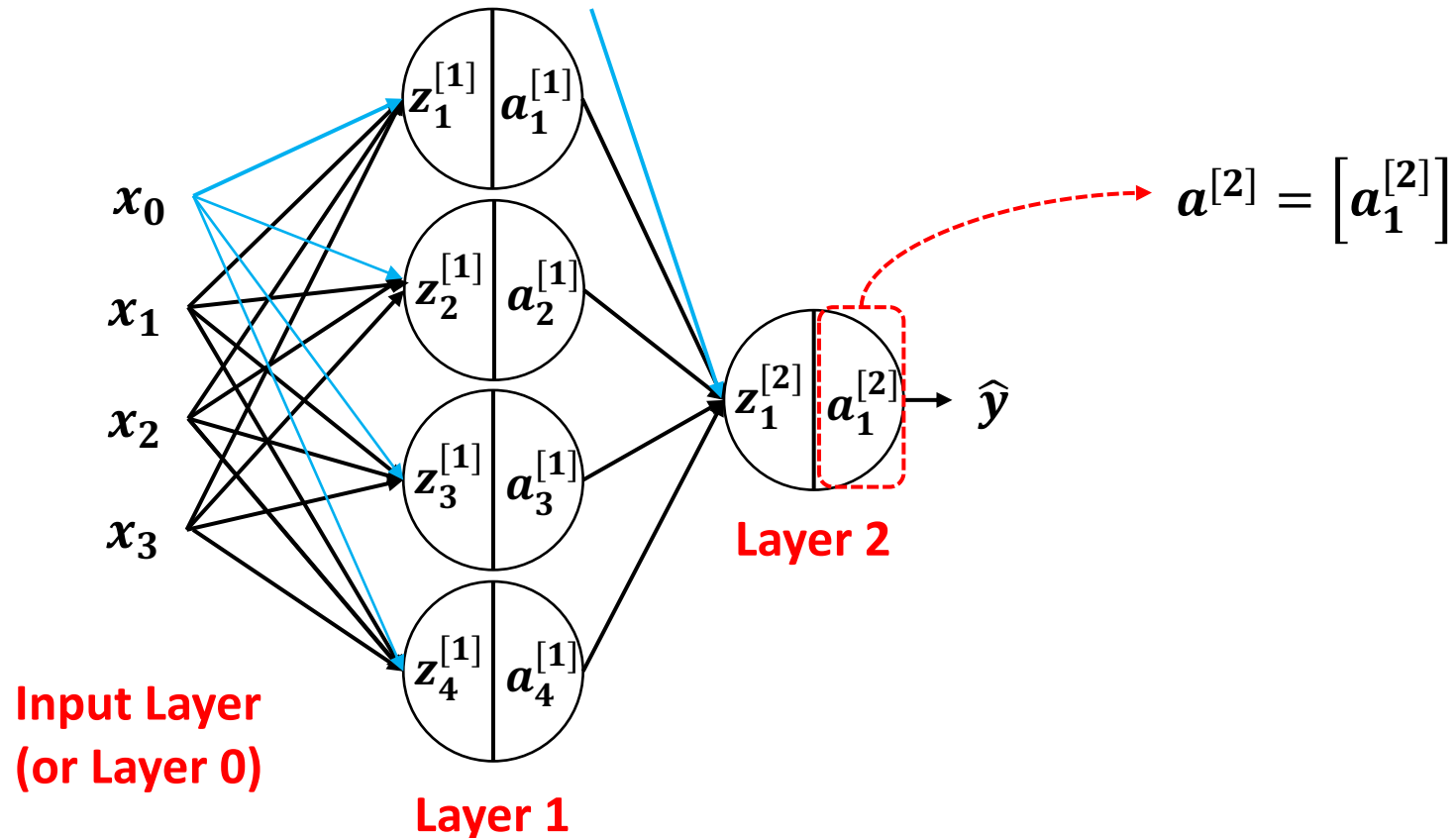
# Vectorizing Input and All Variables

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



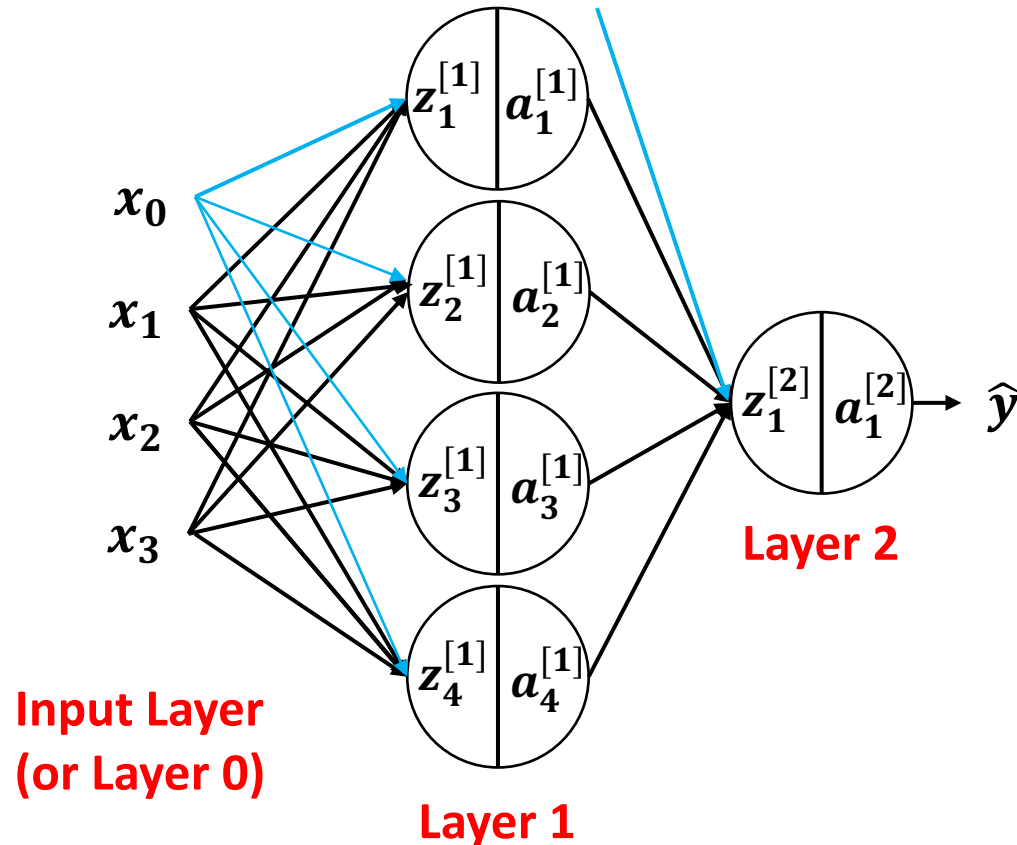
# Vectorizing Input and All Variables

- 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



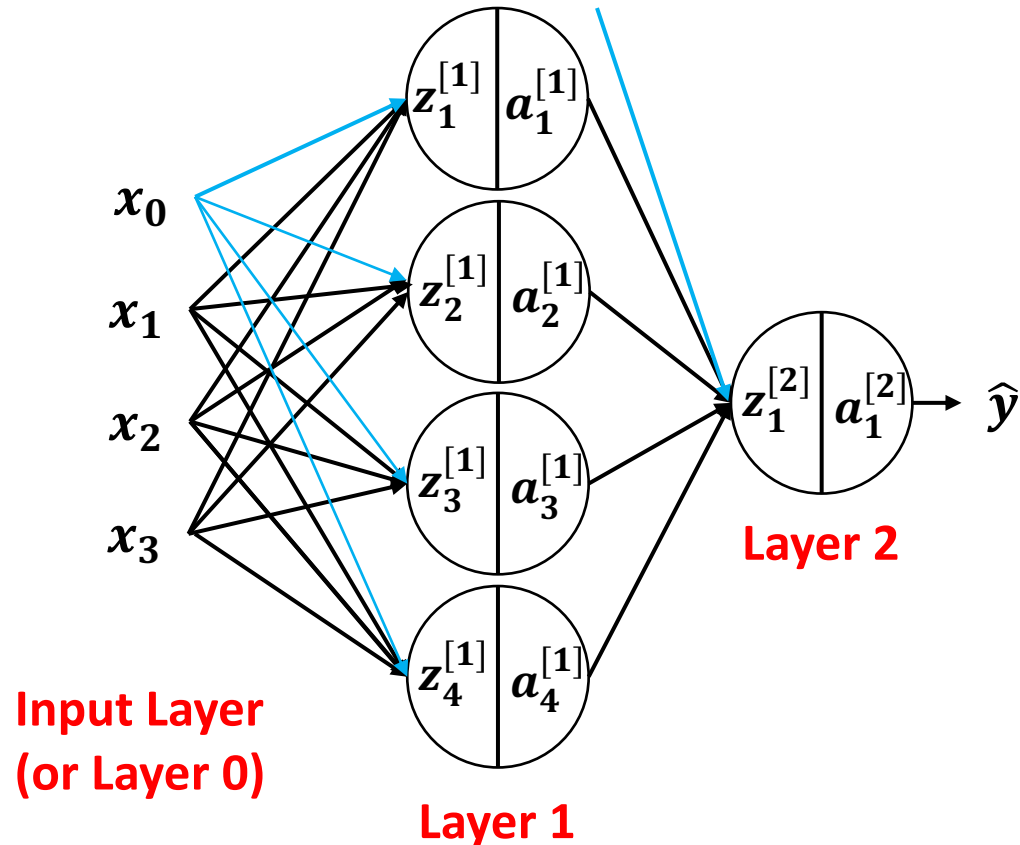
$$\mathbf{z}^{[2]} = \mathbf{w}^{[2]T} \mathbf{x} + \mathbf{b}^{[2]}$$

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

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

# 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



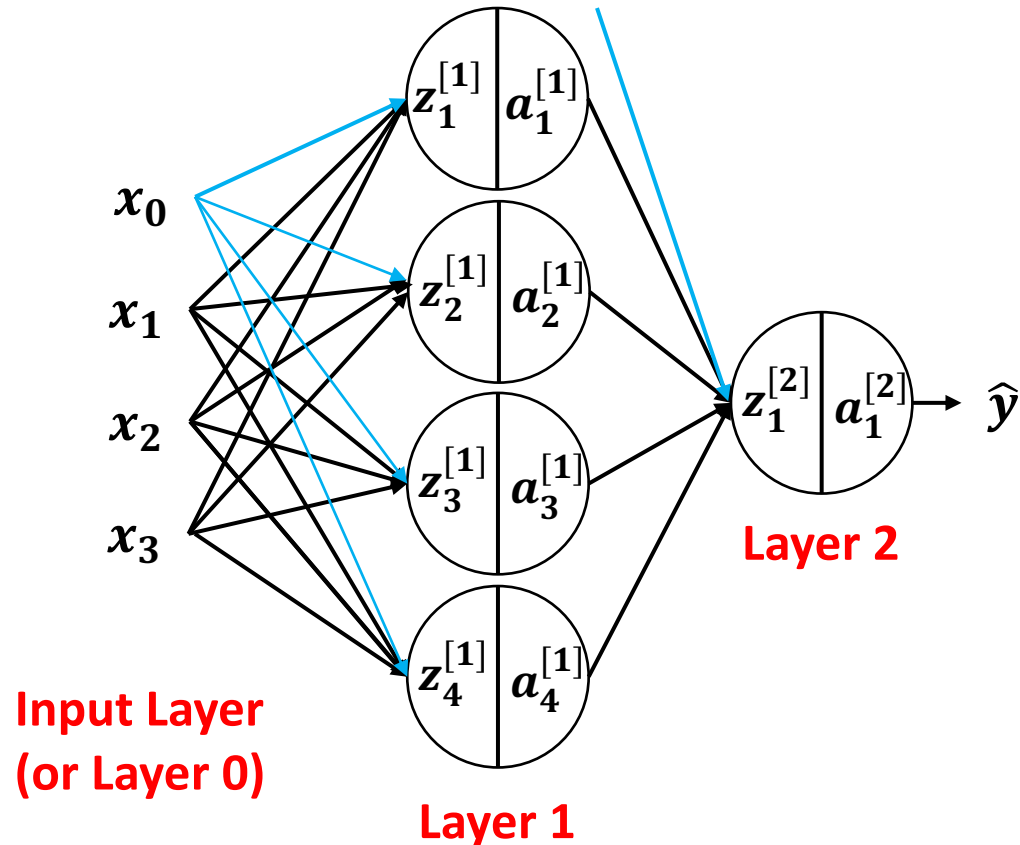
$$\mathbf{z}^{[2]} = \mathbf{w}^{[2]T} \mathbf{x} + \mathbf{b}^{[2]}$$

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

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

# 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



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

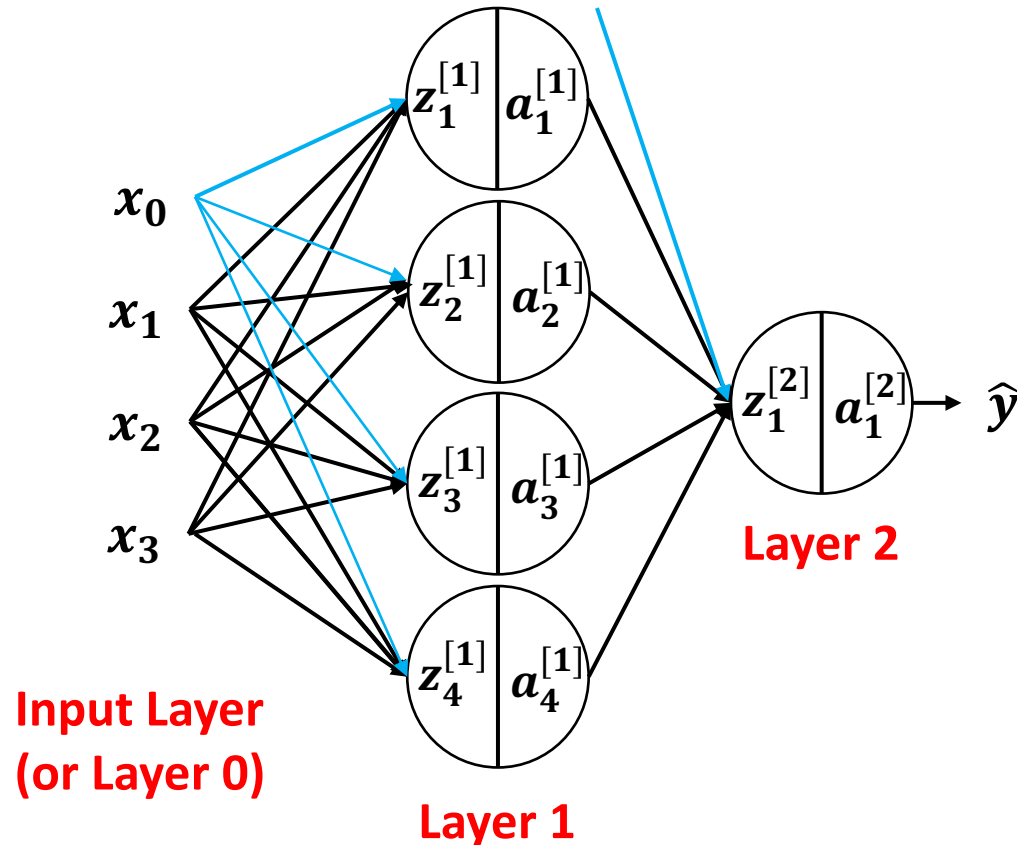
$$= \mathbf{w}^{[2]T} \mathbf{a}^{[1]} + \mathbf{b}^{[2]}$$

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

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

# 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



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

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

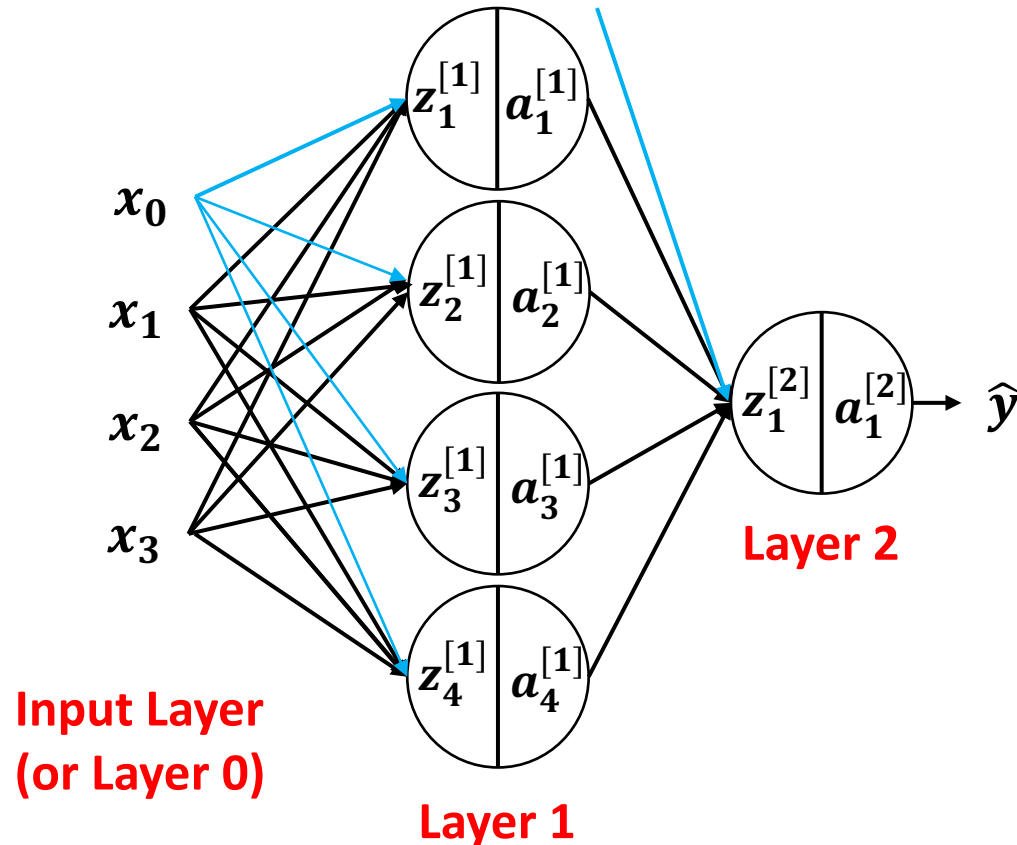
$$= w_{11}^{[2]} a_1^{[1]} + w_{21}^{[2]} a_2^{[1]} + w_{31}^{[2]} a_3^{[1]} + w_{41}^{[2]} a_4^{[1]} + b_1^{[2]}$$

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

$$z_1^{[2]}$$

# 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



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

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

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

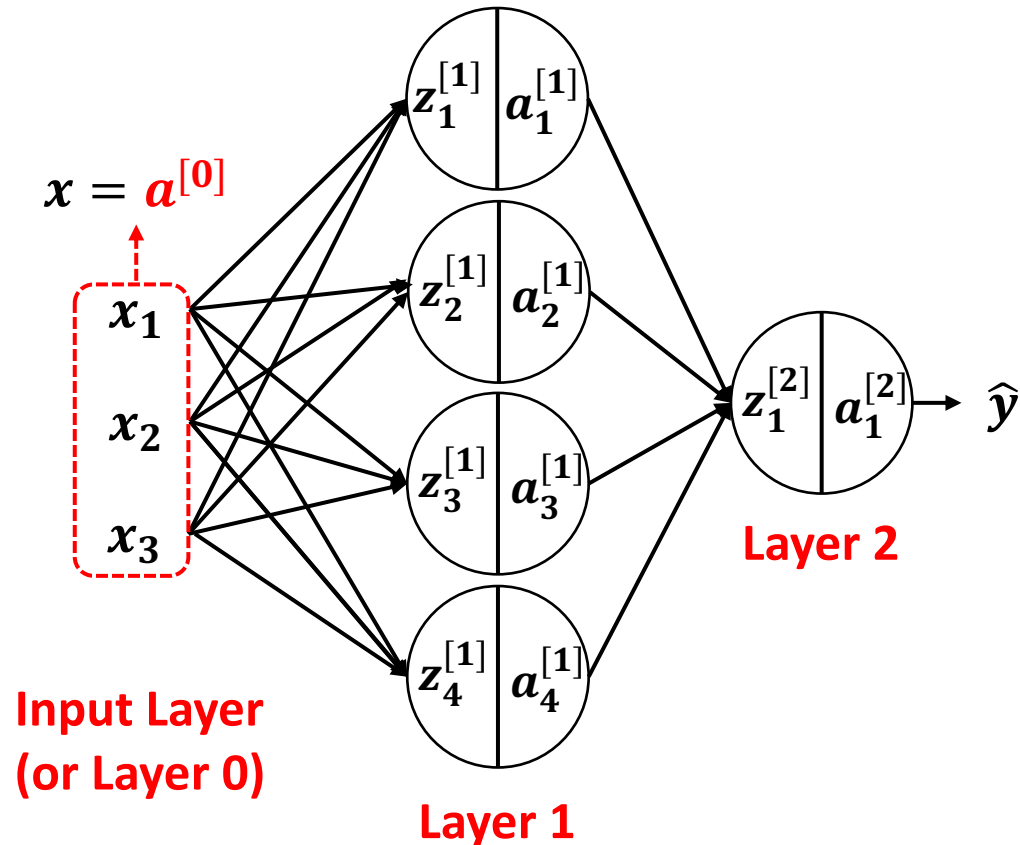
$$\mathbf{a}^{[2]} = \sigma(\mathbf{z}^{[2]})$$

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



# 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



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

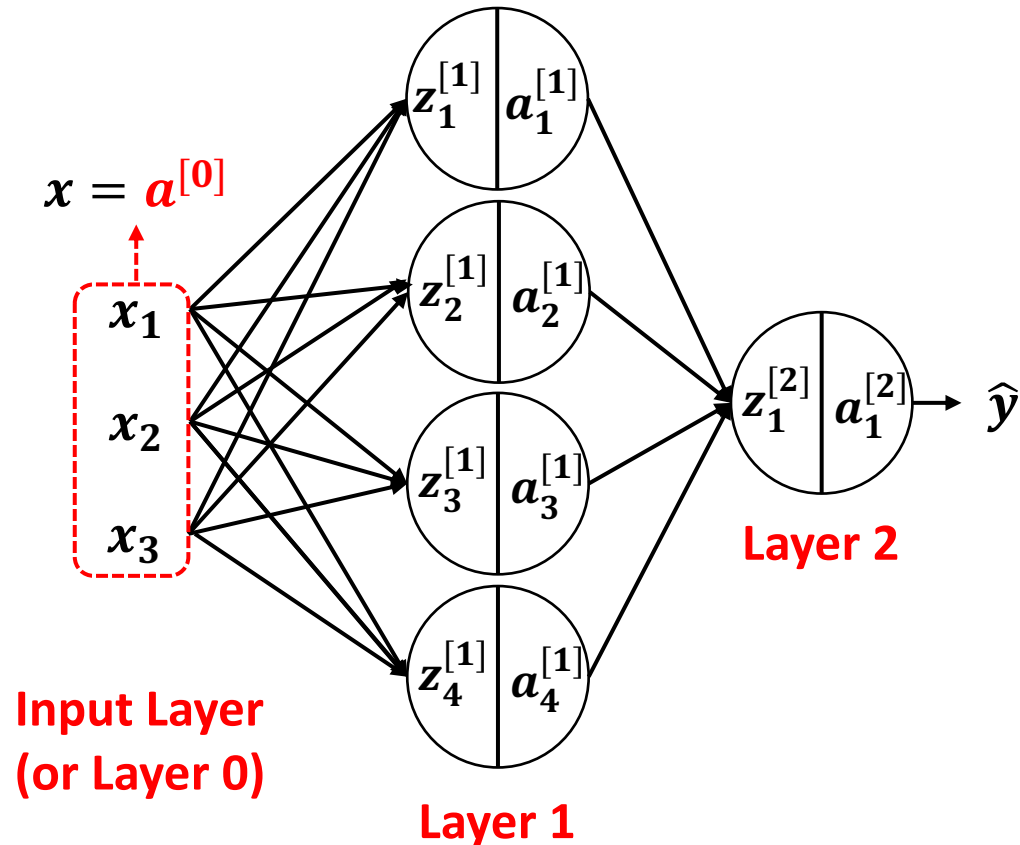
$$a^{[1]} = \sigma(z^{[1]})$$

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

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

# 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



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

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

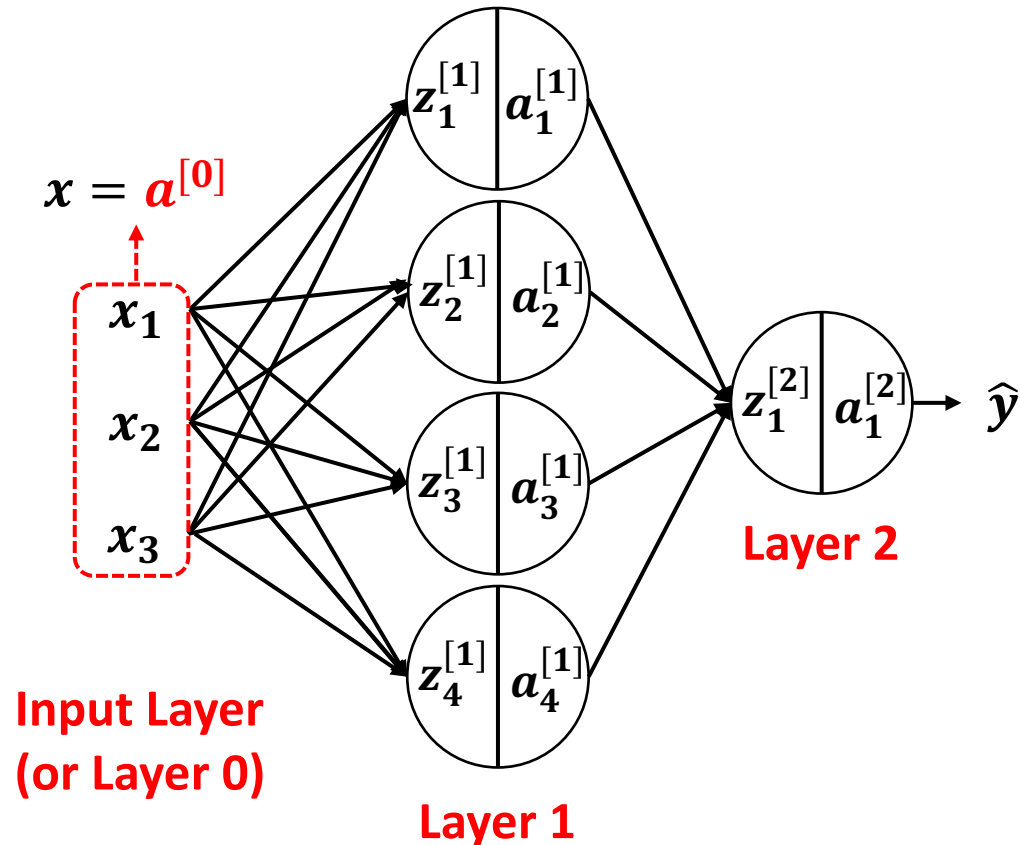
$$\mathbf{z}^{[2]} = \mathbf{w}^{[2]T} \mathbf{a}^{[1]} + \mathbf{b}^{[2]}$$

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

But, this assumes only 1 training example!

# 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



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

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

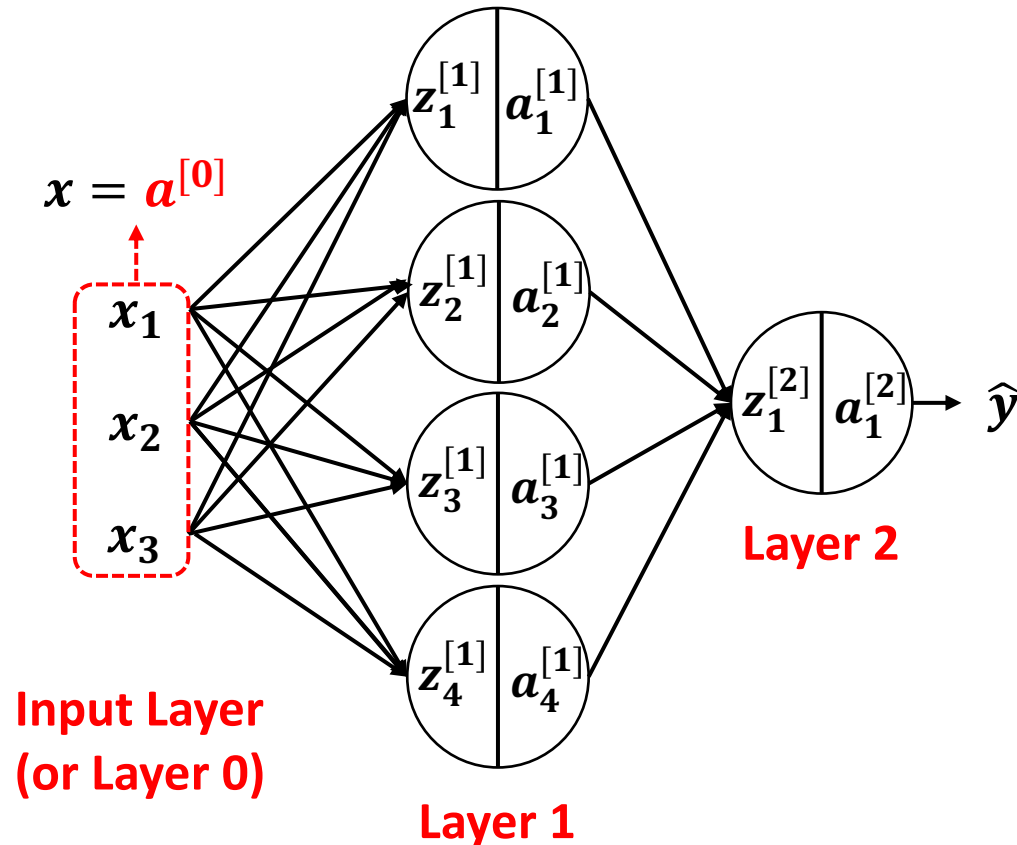
$$\mathbf{z}^{[2]} = \mathbf{w}^{[2]T} \mathbf{a}^{[1]} + \mathbf{b}^{[2]}$$

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

How can we account for all the training examples?

# 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



*for*  $i = 1$  to  $n$ : Assuming  $n$  examples

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

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

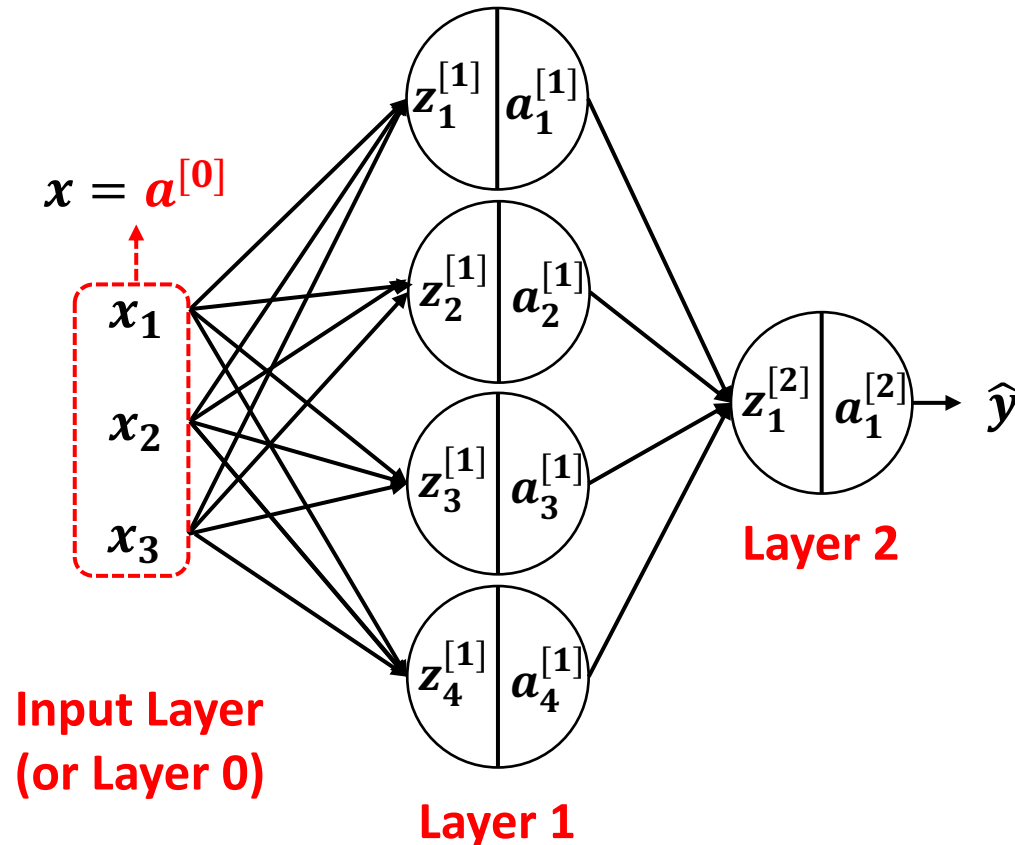
$$z^{[2]}(i) = w^{[2]T}(i) a^{[1]}(i) + b^{[2]}$$

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

Refers to the  $i^{th}$  example in the training dataset

# 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



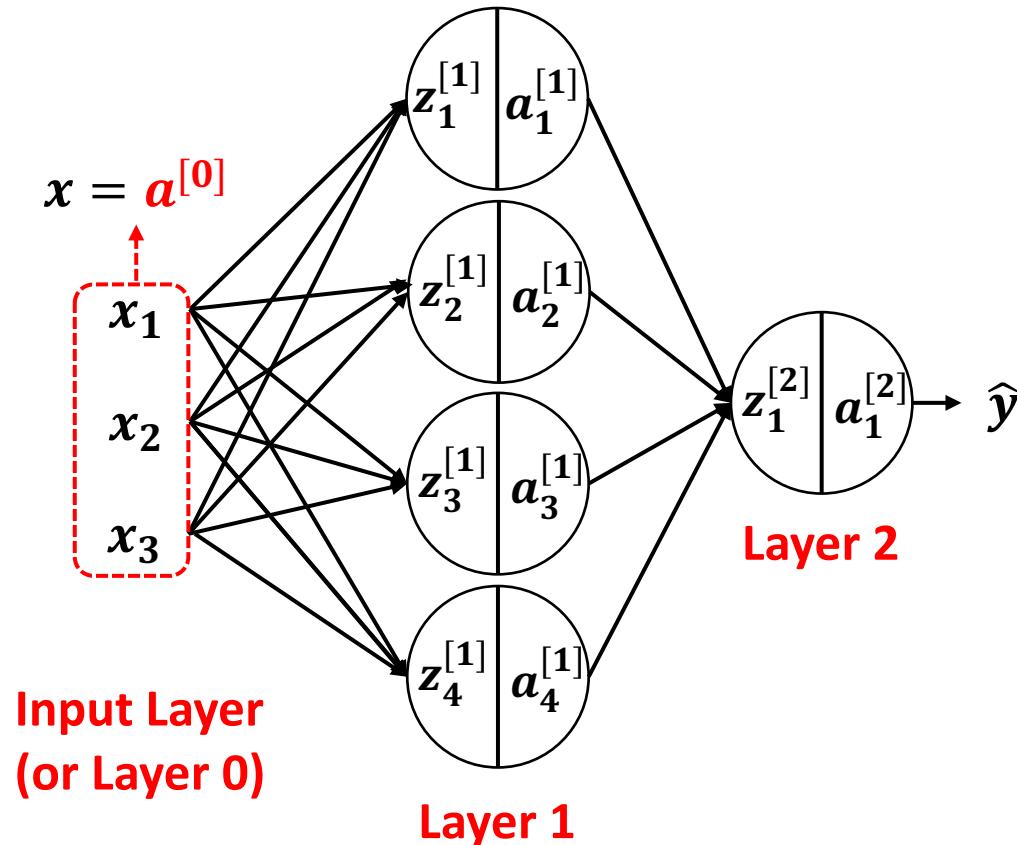
Assuming  $n$  examples  
for  $i = 1$  to  $n$ :

But, loops in general slow down programs; hence, it is better to further **vectorize** the implementation in order to avoid any loop, whenever possible

Refers to the  $i^{\text{th}}$  example in the training dataset

# 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



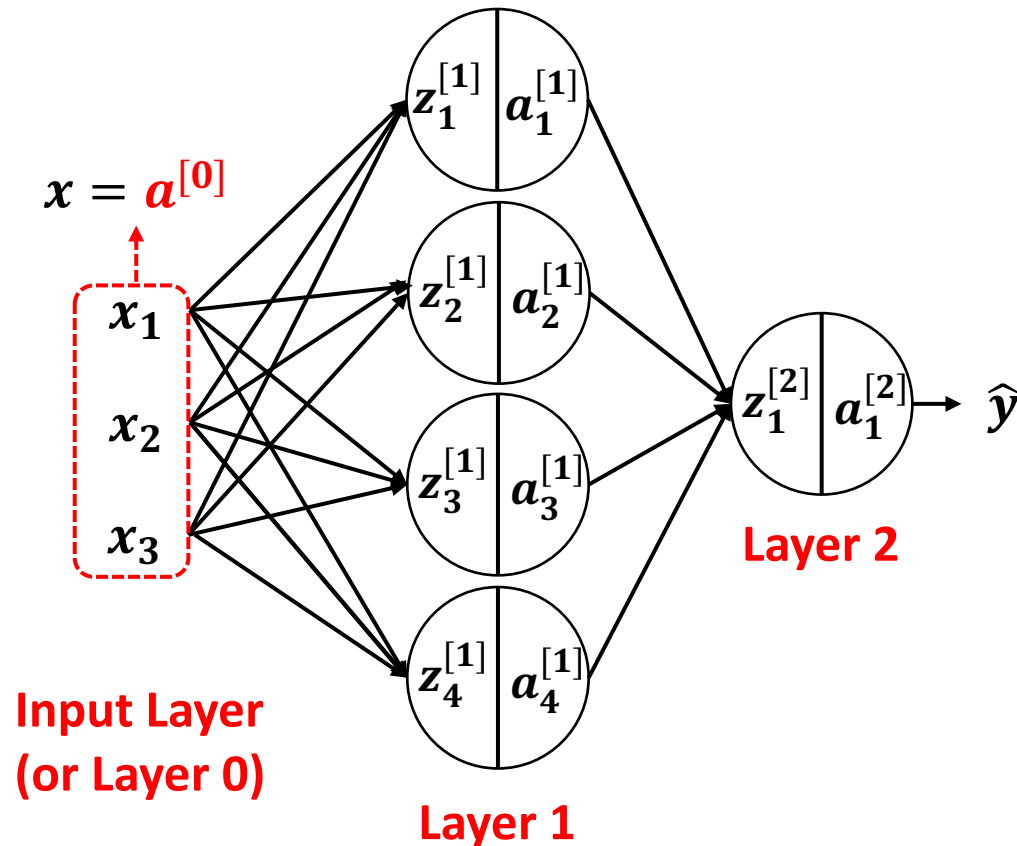
Assuming  $n$  examples  
for  $i = 1$  to  $n$ :

To this end, we can simply stack all  $x$  vectors (or  $a^{[0]}$  vectors),  $z$  vectors, and  $a$  vectors in different matrices of every layer!

Refers to the  $i^{th}$  example in the training dataset

# 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



$$X = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(n)} \\ x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(n)} \\ x_3^{(1)} & x_3^{(2)} & \dots & x_3^{(n)} \end{bmatrix}$$

A red dashed box highlights the first column of the matrix  $X$ , which contains the elements  $x_1^{(1)}, x_2^{(1)}, x_3^{(1)}$ . A red arrow points from this box to the text below.

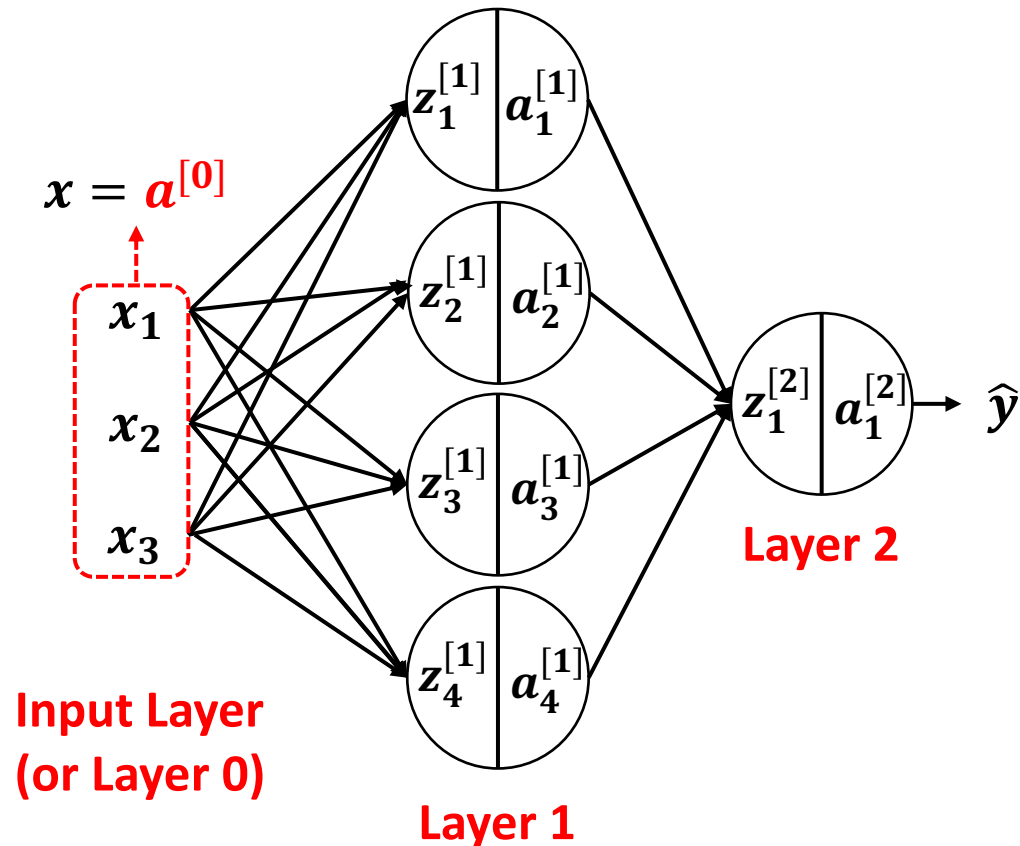
A green dashed arrow points from the matrix  $X$  to the text on the right.

This vector represents the *first* example in the training dataset

Assuming  $n$  examples

# 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



$$X = \begin{bmatrix} x_1^{(1)} & \boxed{x_1^{(2)}} & \dots & x_1^{(n)} \\ x_2^{(1)} & \boxed{x_2^{(2)}} & \dots & x_2^{(n)} \\ x_3^{(1)} & \boxed{x_3^{(2)}} & \dots & x_3^{(n)} \end{bmatrix}$$

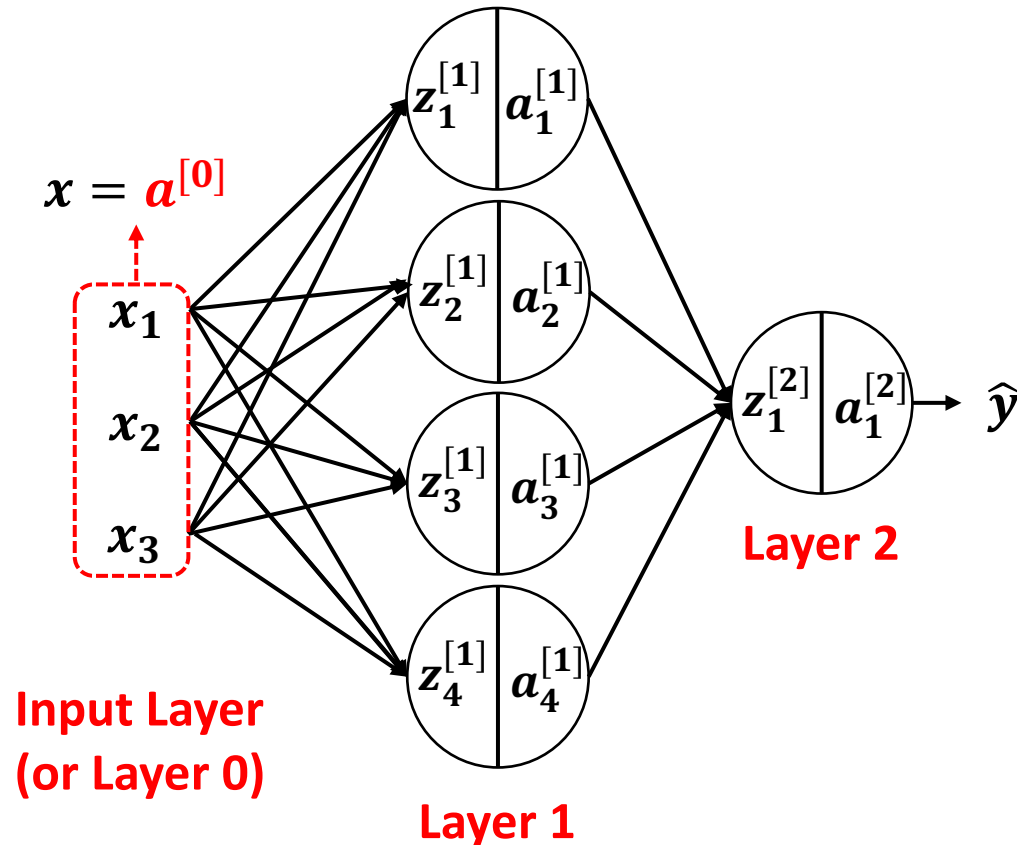
Assuming  $n$  examples

This vector represents the *second* example in the training dataset



# 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



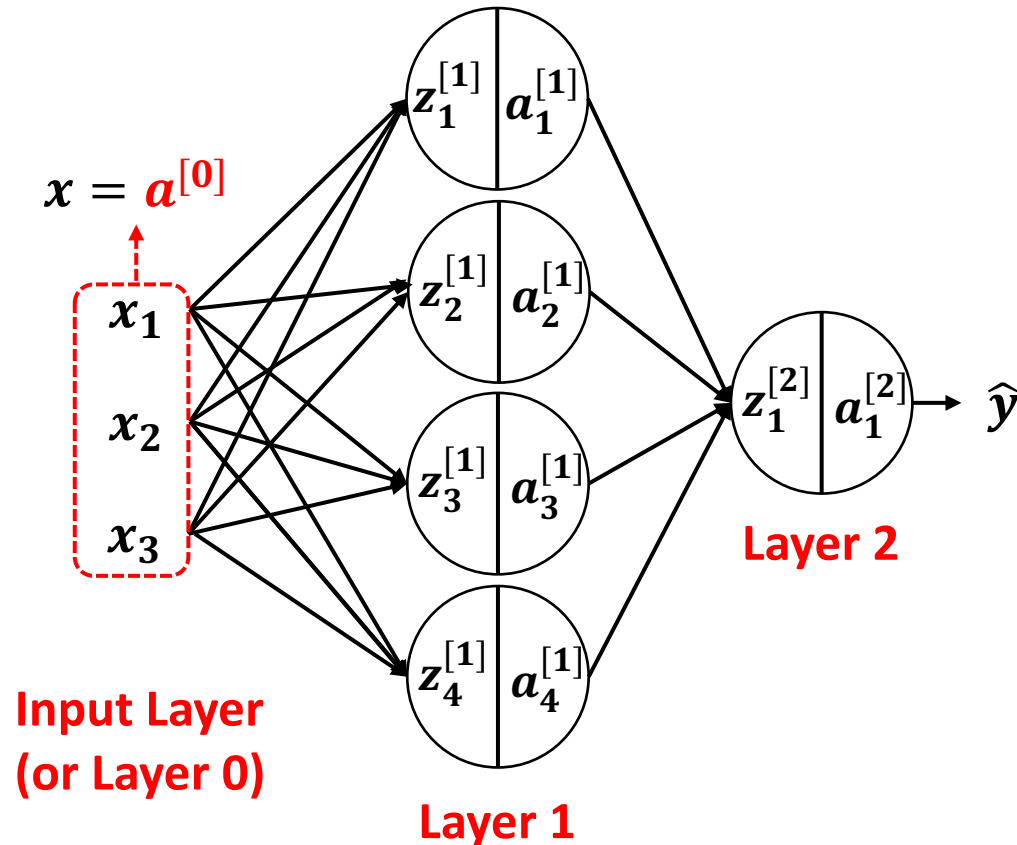
$$X = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(n)} \\ x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(n)} \\ x_3^{(1)} & x_3^{(2)} & \dots & x_3^{(n)} \end{bmatrix}$$

Assuming  $n$  examples

This vector represents the  $n^{th}$  example in the training dataset

# 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

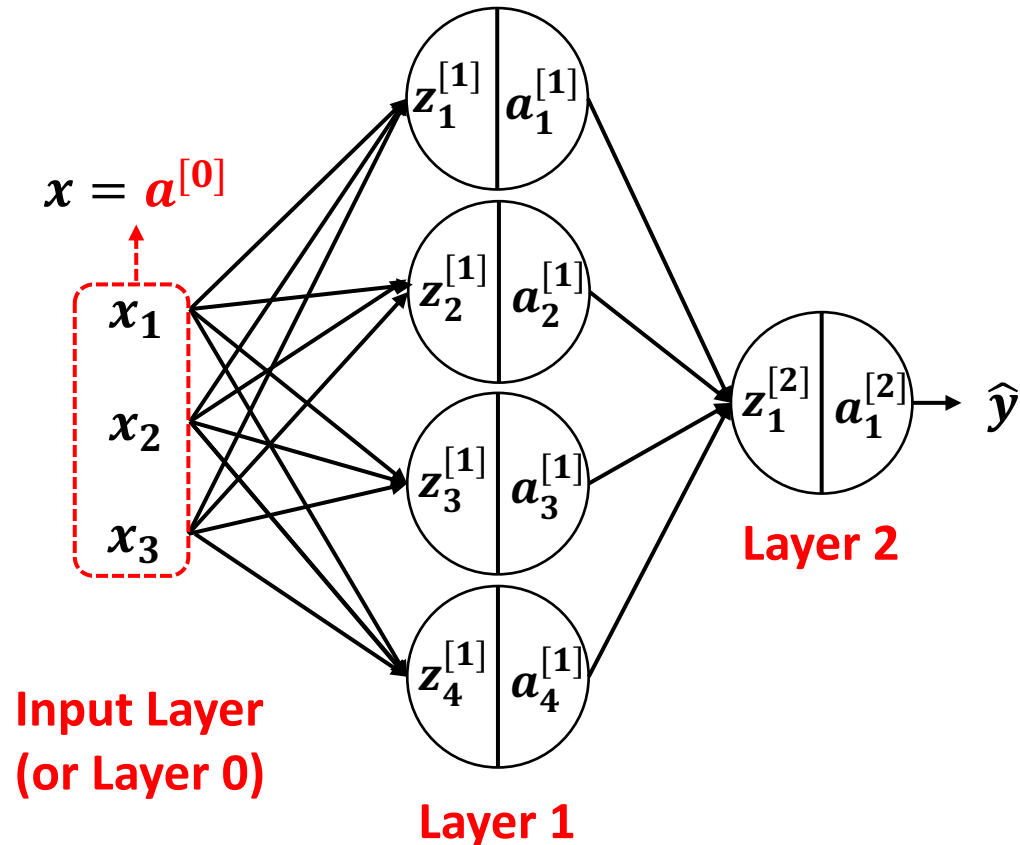


$$X = \begin{bmatrix} a_1^{[0](1)} & a_1^{[0](2)} & \dots & a_1^{[0](n)} \\ a_2^{[0](1)} & a_2^{[0](2)} & \dots & a_2^{[0](n)} \\ a_3^{[0](1)} & a_3^{[0](2)} & \dots & a_3^{[0](n)} \end{bmatrix}$$

If we denote  $x$  as  $a^{[0]}$

# 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

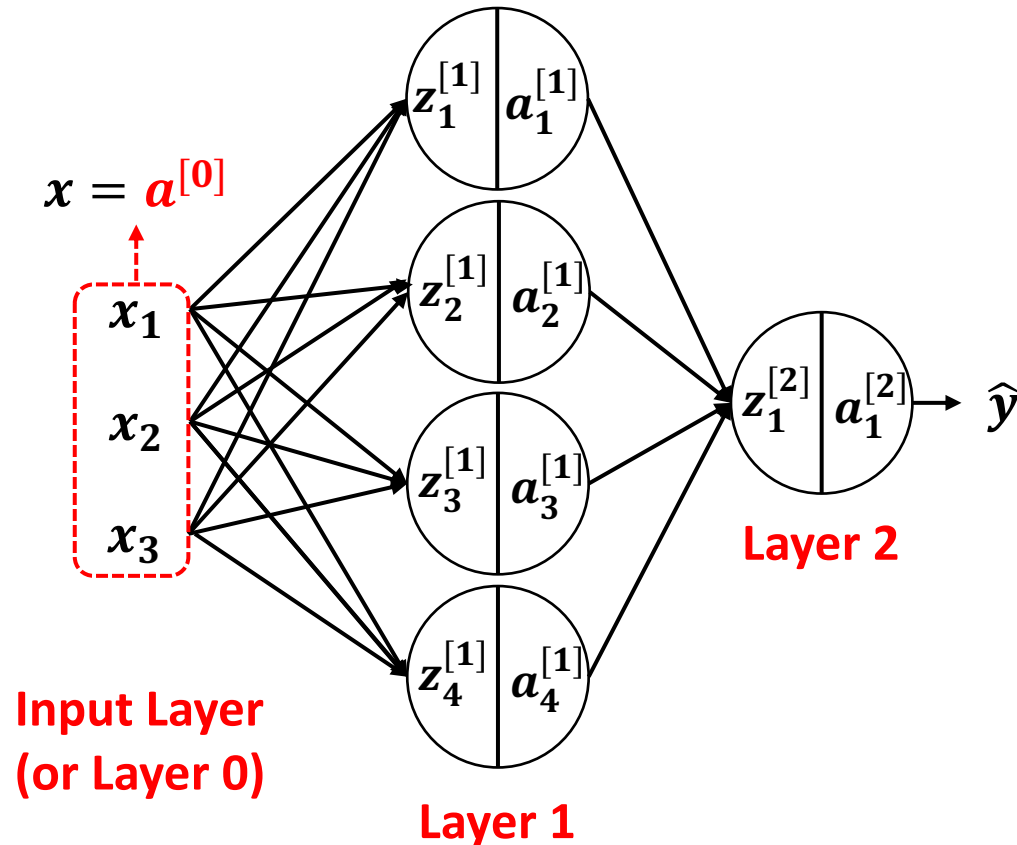


$$X = A^{[0]} = \begin{bmatrix} a_1^{[0](1)} & a_1^{[0](2)} & \dots & a_1^{[0](n)} \\ a_2^{[0](1)} & a_2^{[0](2)} & \dots & a_2^{[0](n)} \\ a_3^{[0](1)} & a_3^{[0](2)} & \dots & a_3^{[0](n)} \end{bmatrix}$$

If we denote  $x$  as  $a^{[0]}$

# Vectorizing Input and All Variables

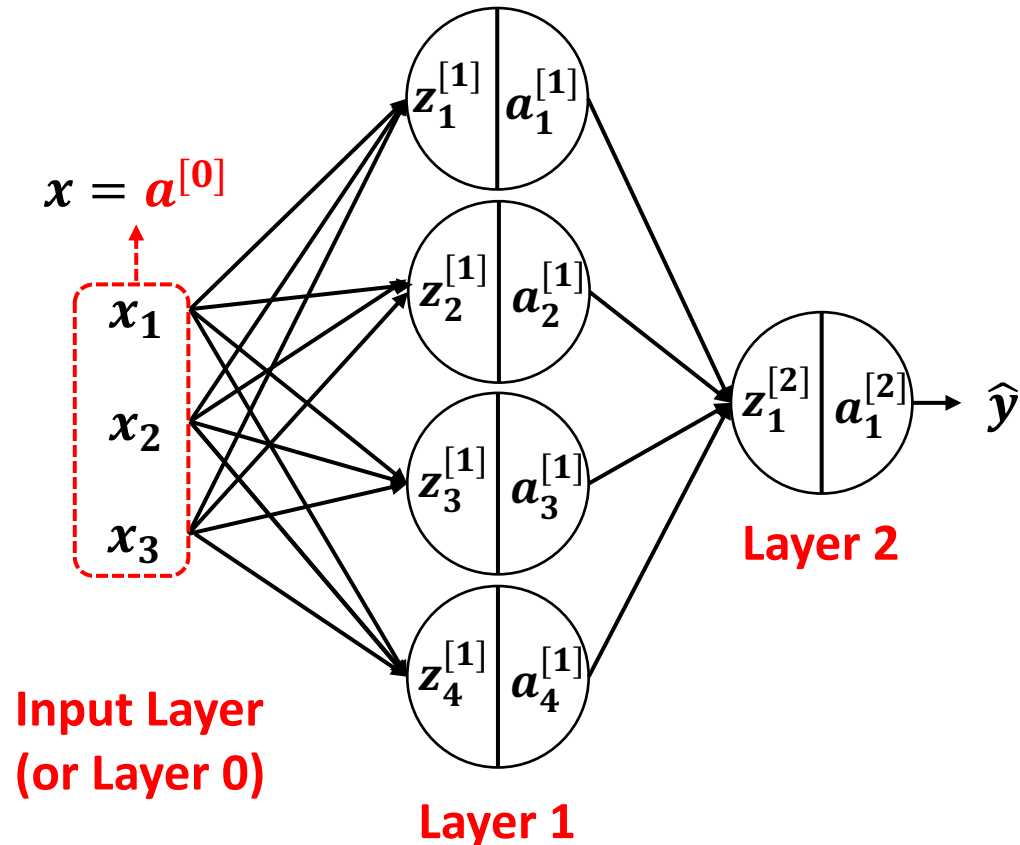
- To help develop an efficient learning algorithm, let us **vectorize** (represent in vectors & matrices) the input and the variables involved



$$Z^{[1]} = \begin{bmatrix} z_1^{[1]}(1) & z_1^{[1]}(2) & \dots & z_1^{[1]}(n) \\ z_2^{[1]}(1) & z_2^{[1]}(2) & \dots & z_2^{[1]}(n) \\ z_3^{[1]}(1) & z_3^{[1]}(2) & \dots & z_3^{[1]}(n) \\ z_4^{[1]}(1) & z_4^{[1]}(2) & \dots & z_4^{[1]}(n) \end{bmatrix}$$

# Vectorizing Input and All Variables

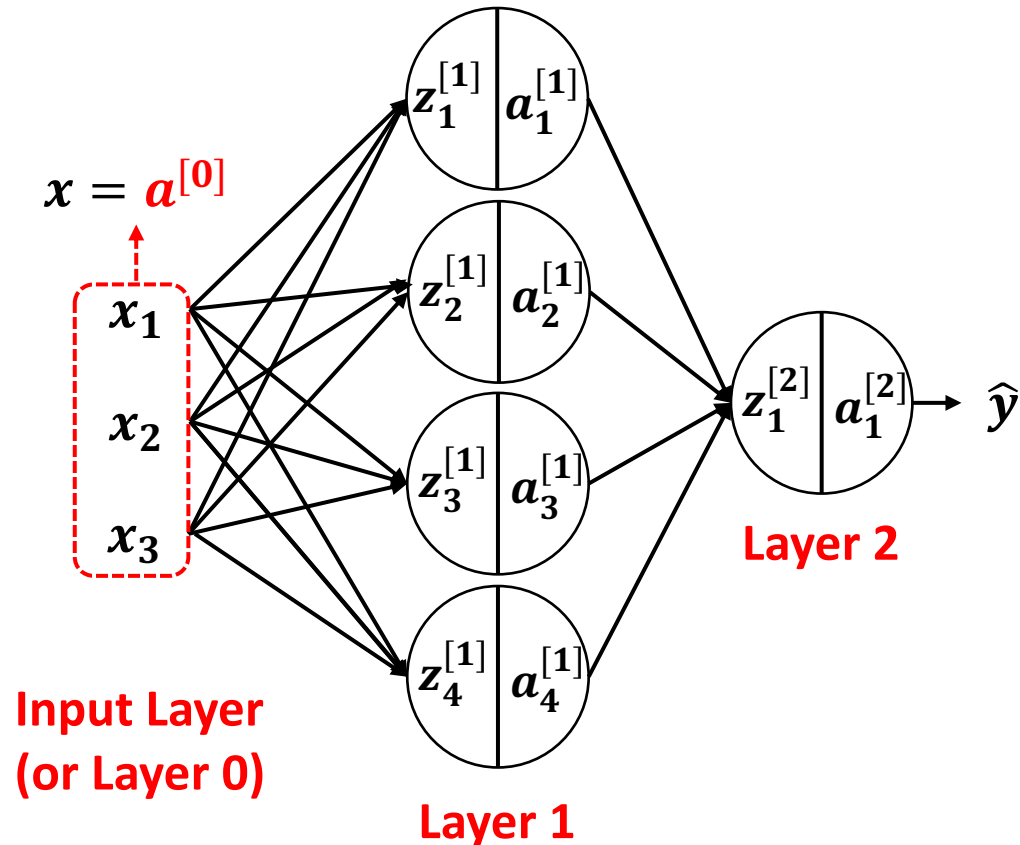
- To help develop an efficient learning algorithm, let us **vectorize** (represent in vectors & matrices) the input and the variables involved



$$A^{[1]} = \begin{bmatrix} a_1^{[1](1)} & a_1^{[1](2)} & \dots & a_1^{[1](n)} \\ a_2^{[1](1)} & a_2^{[1](2)} & \dots & a_2^{[1](n)} \\ a_3^{[1](1)} & a_3^{[1](2)} & \dots & a_3^{[1](n)} \\ a_4^{[1](1)} & a_4^{[1](2)} & \dots & a_4^{[1](n)} \end{bmatrix}$$

# Vectorizing Input and All Variables

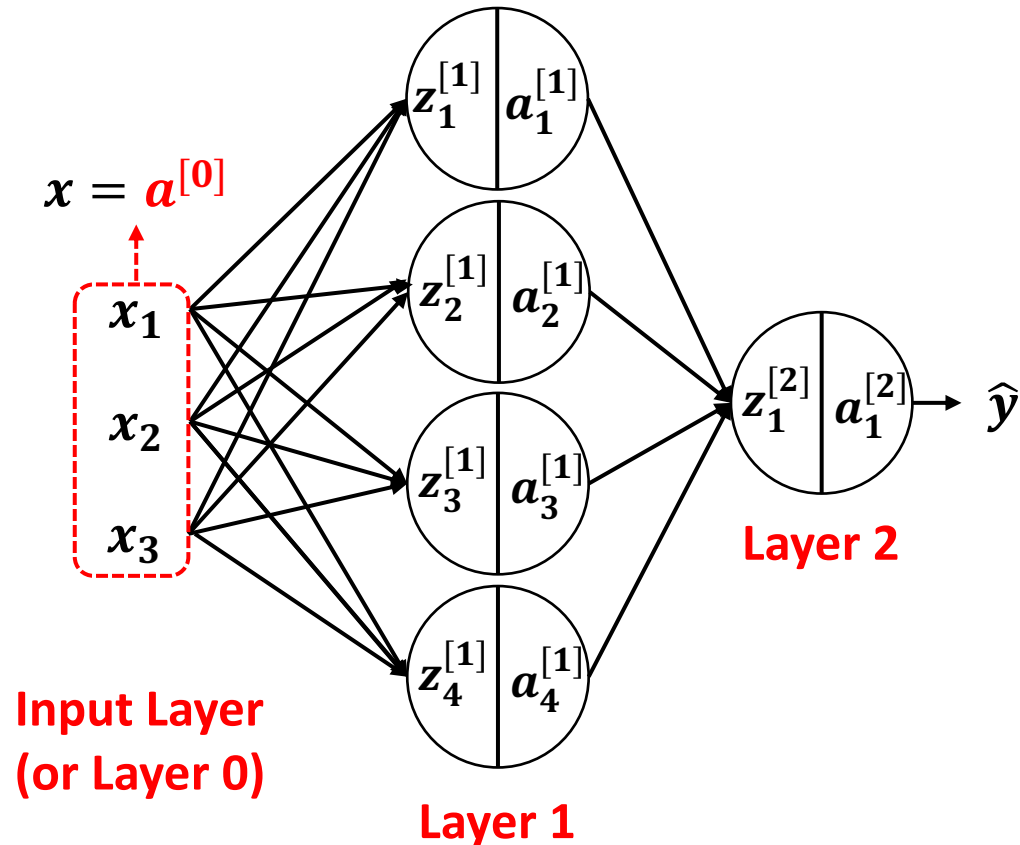
- To help develop an efficient learning algorithm, let us **vectorize** (represent in vectors & matrices) the input and the variables involved



$$Z^{[2]} = \begin{bmatrix} z_1^{[2]}(1) & z_1^{[2]}(2) & \dots & z_1^{[2]}(n) \end{bmatrix}$$

# Vectorizing Input and All Variables

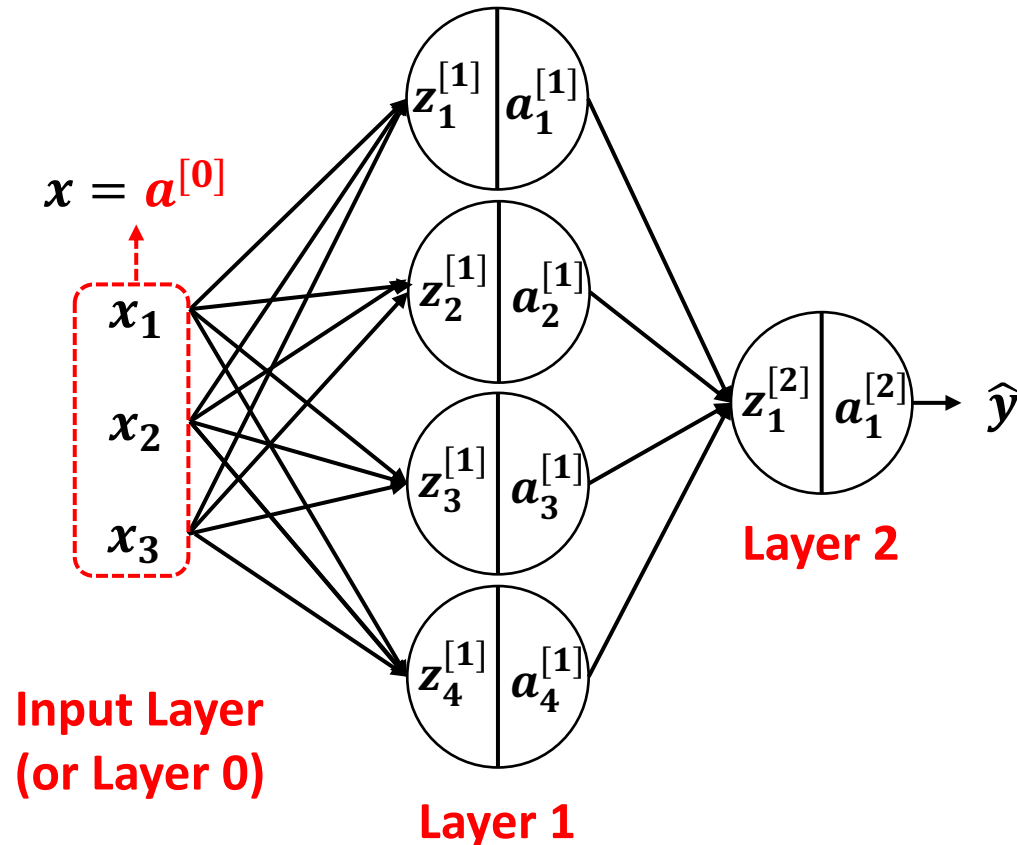
- To help develop an efficient learning algorithm, let us **vectorize** (represent in vectors & matrices) the input and the variables involved



$$A^{[2]} = \begin{bmatrix} a_1^{[2](1)} & a_1^{[2](2)} & \dots & a_1^{[2](n)} \end{bmatrix}$$

# 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



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

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

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

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

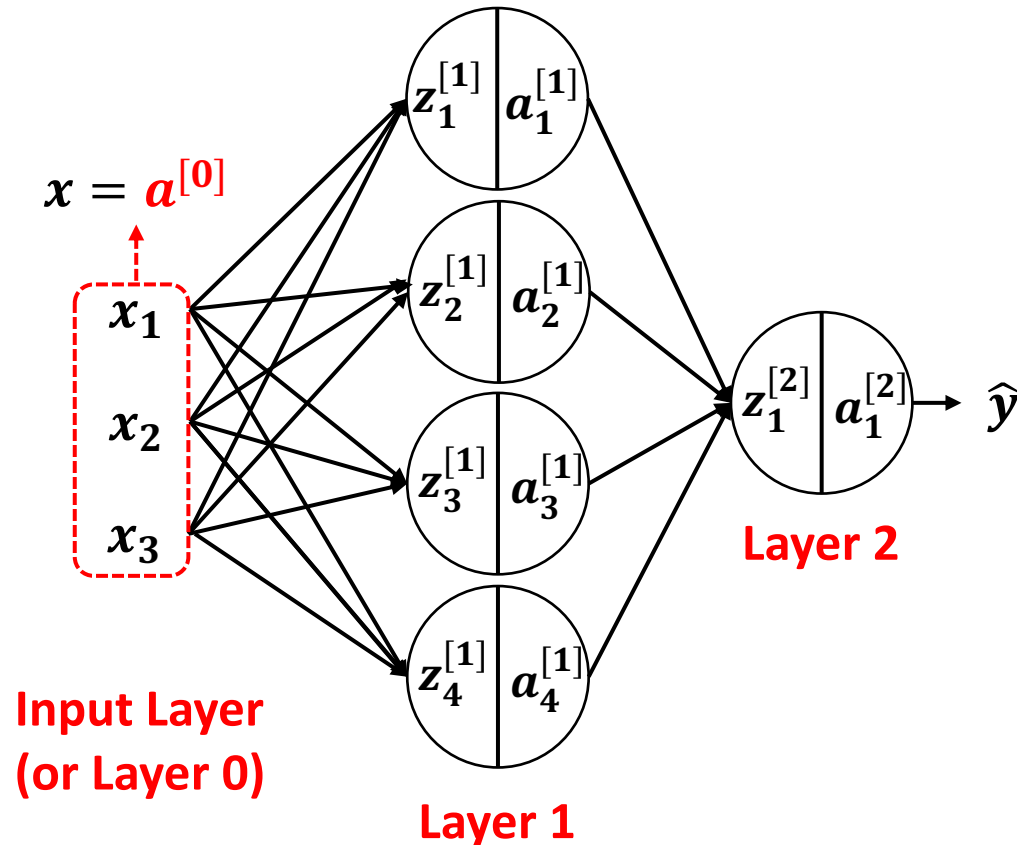
$$a^{[2]}(i) = \sigma(z^{[2]}(i))$$

**Before Vectorization**



# 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



**No Explicit Loop!**

$$Z^{[1]} = w^{[1]T} A^{[0]} + b^{[1]}$$

$$A^{[1]} = \sigma(Z^{[1]})$$

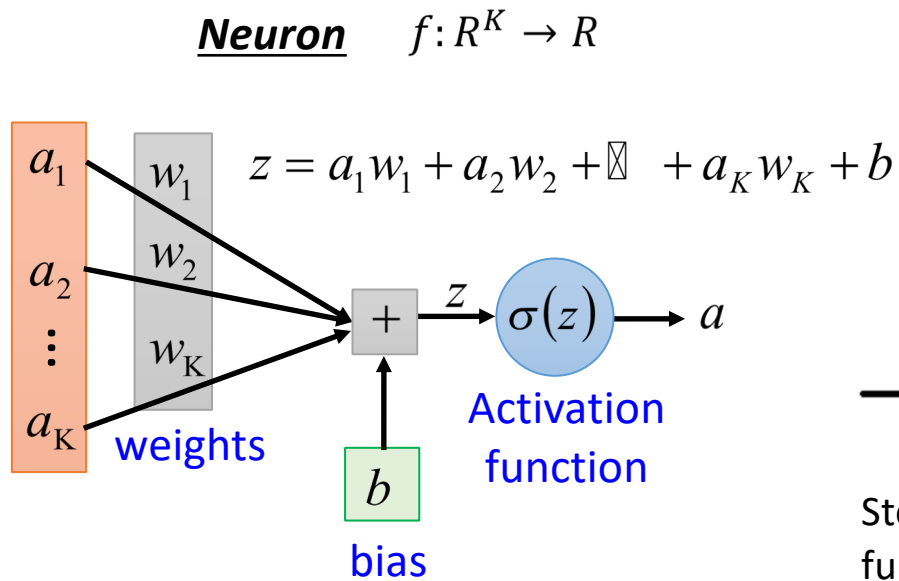
$$Z^{[2]} = w^{[2]T} A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(Z^{[2]})$$

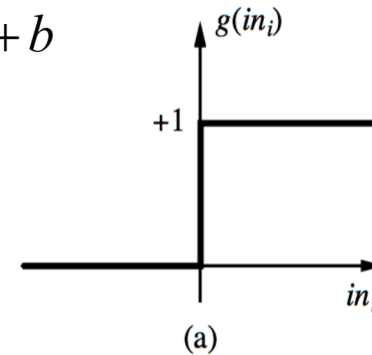
**After Vectorization**

# Epitome: A Neuron

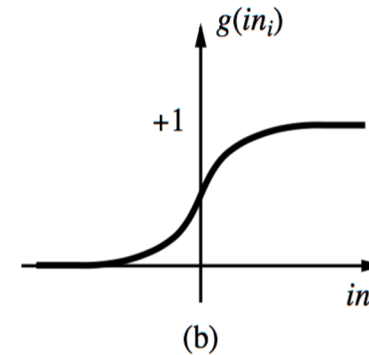
A neuron is a computational unit in the neural network that exchanges messages with each other.



Possible activation functions:



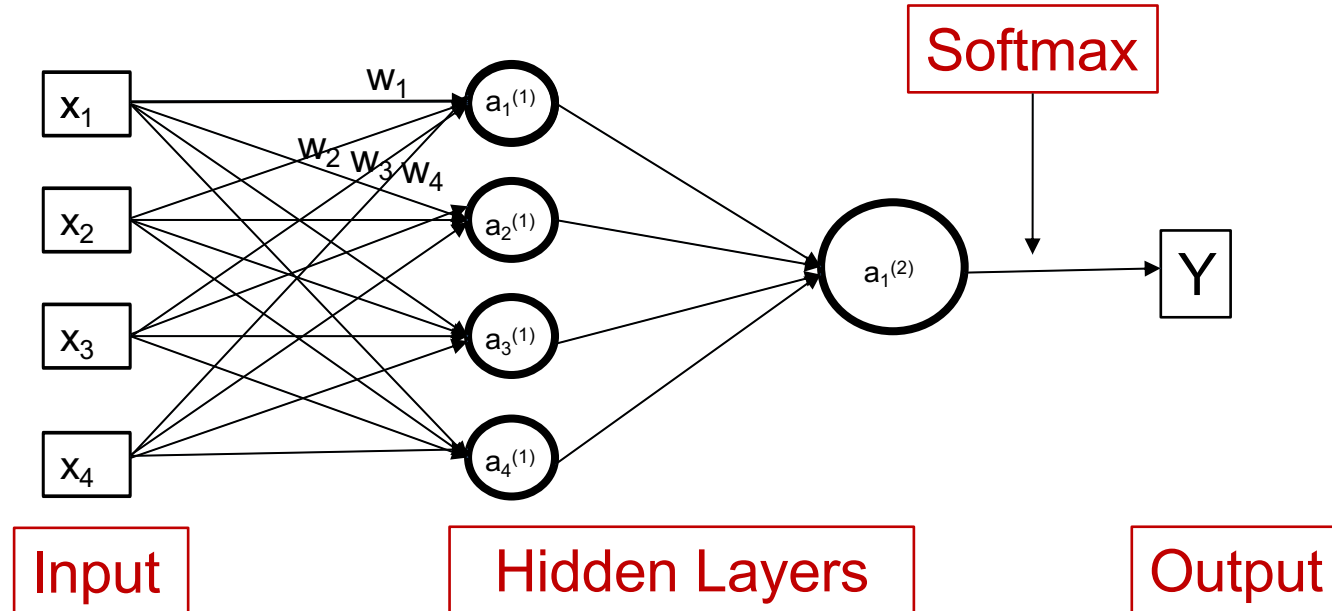
Step  
function/threshold  
function



Sigmoid function (a.k.a,  
logistic function)

# Epitome: A simple neural network

Ignoring bias...



$$a_1^{(1)} = f(w_1 * x_1 + w_2 * x_2 + w_3 * x_3 + w_4 * x_4)$$

Number of  
parameters:

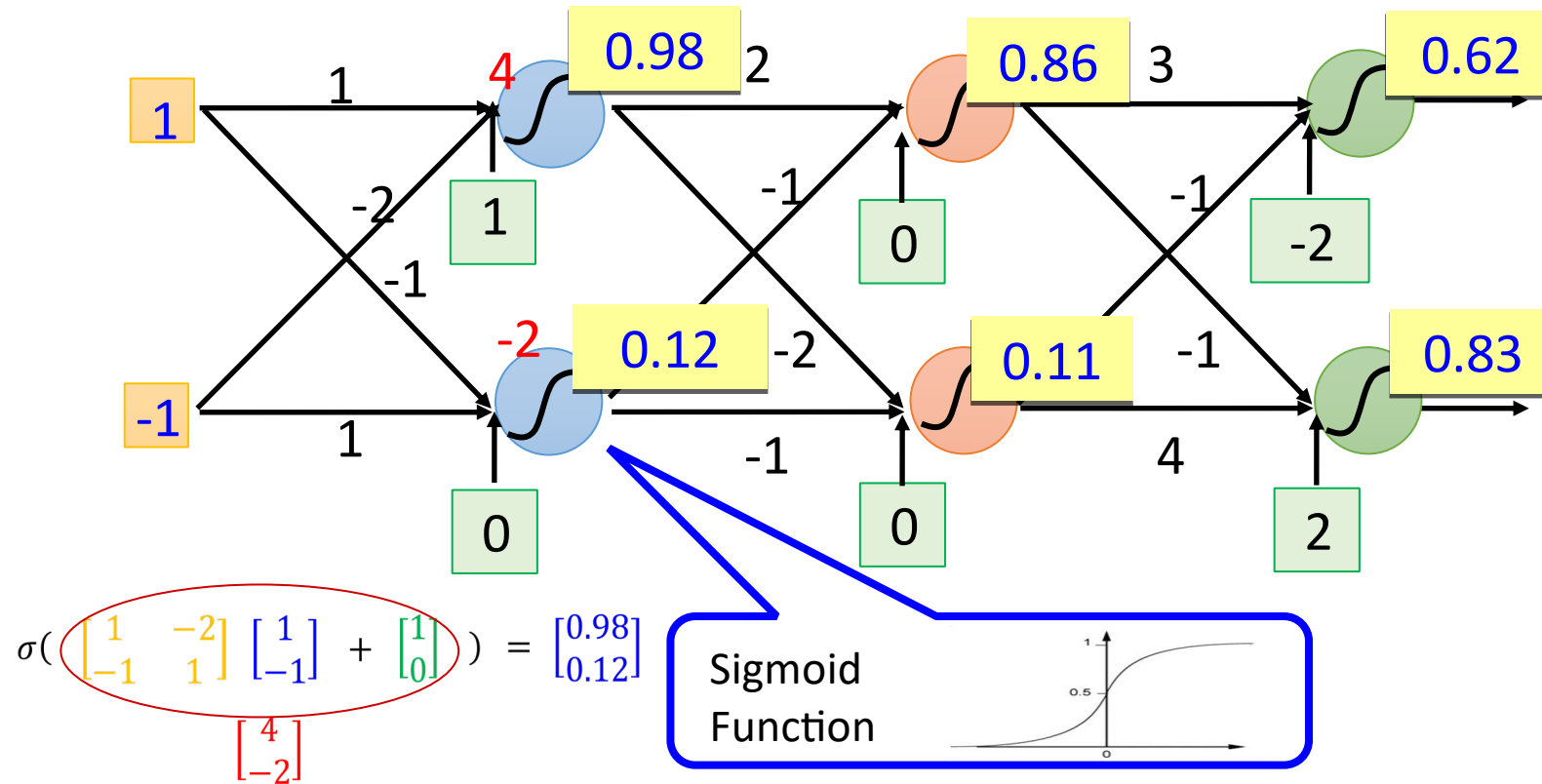
$$4 * 4 + 4 + 1$$

$f()$  is activation function: Relu or sigmoid

*Relu*:  $\max(0, x)$

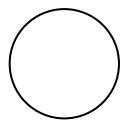
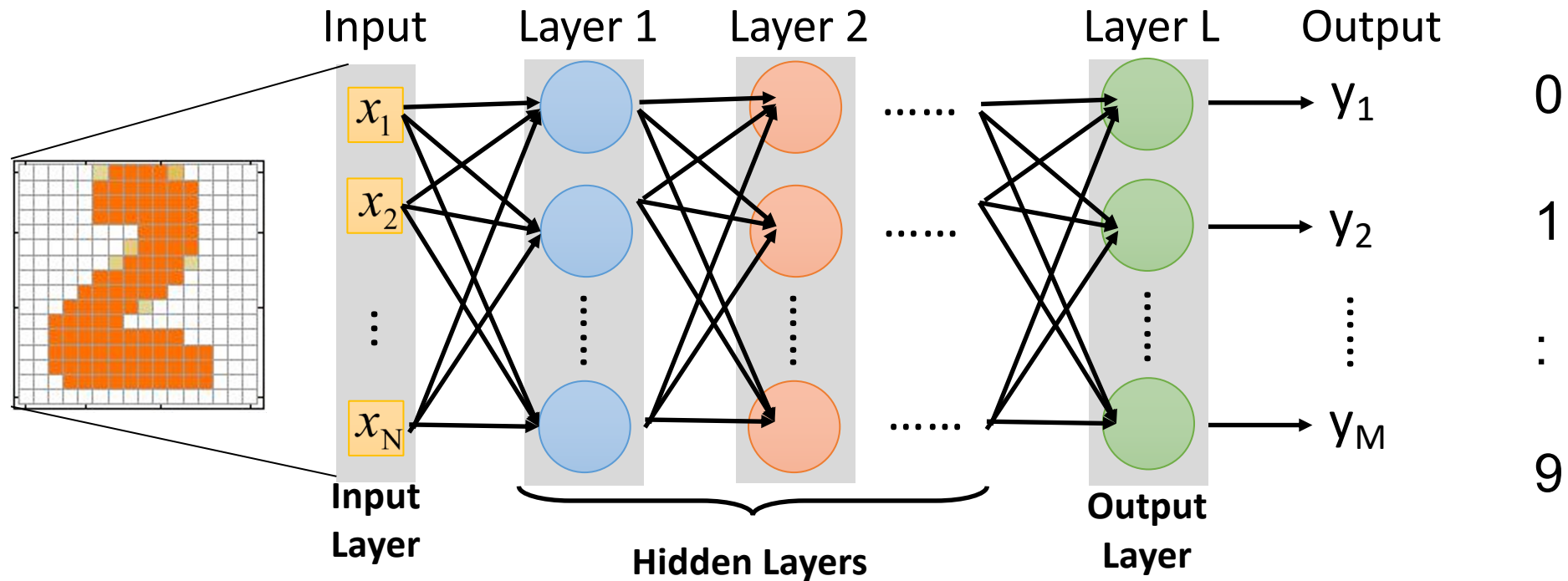
$$a_1^{(1)} = \max(0, w_1 * x_1 + w_2 * x_2 + w_3 * x_3 + w_4 * x_4)$$

# Epitome: Example of Neural Network



$$f: \mathbb{R}^2 \rightarrow \mathbb{R}^2 \quad f\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix}\right) = \begin{bmatrix} 0.62 \\ 0.83 \end{bmatrix} \quad f\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 0.51 \\ 0.85 \end{bmatrix}$$

# Epitome: A case study



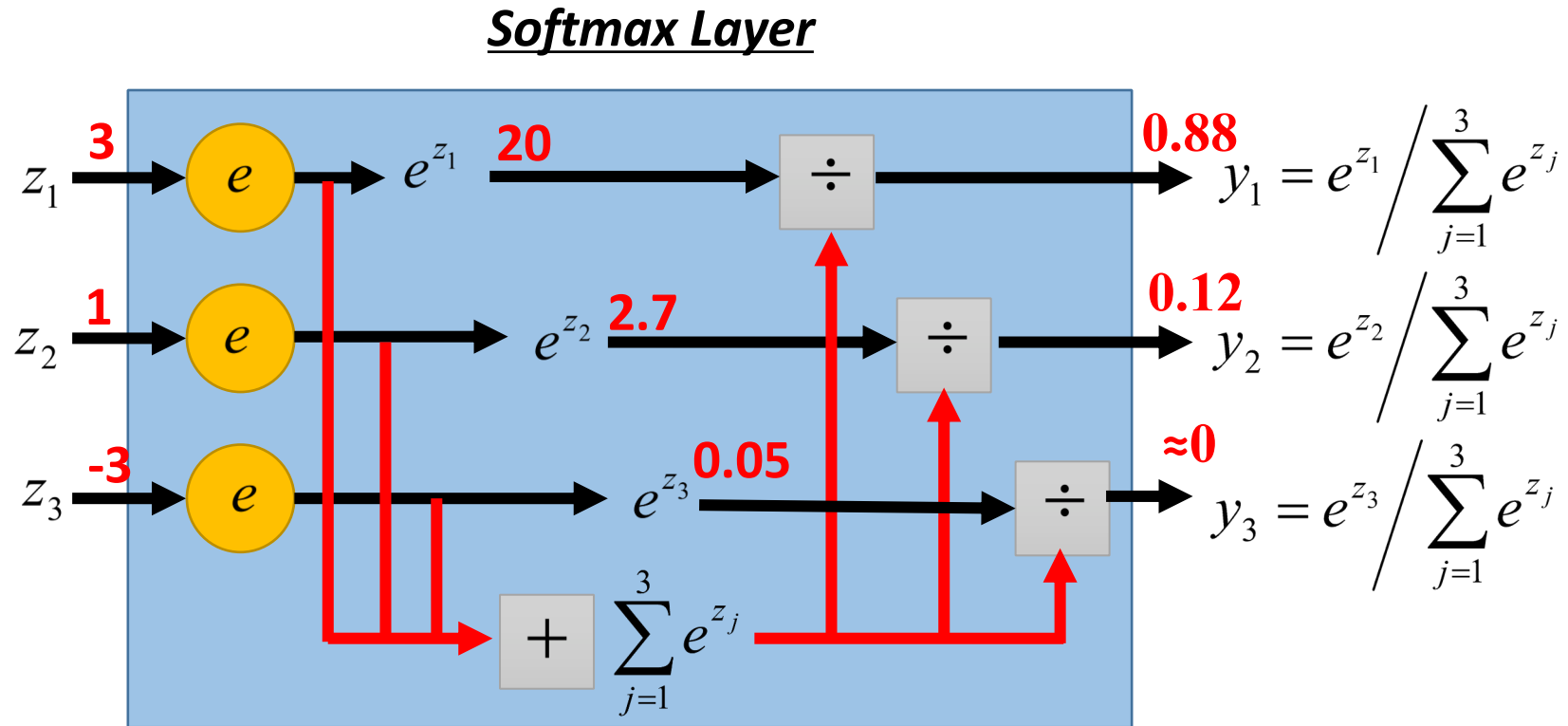
: denotes neuron

$$f: R^{256} \rightarrow R^{10}$$

Deep learning (DL) means many hidden layers.  
DL represents the function  $f$  by neural network

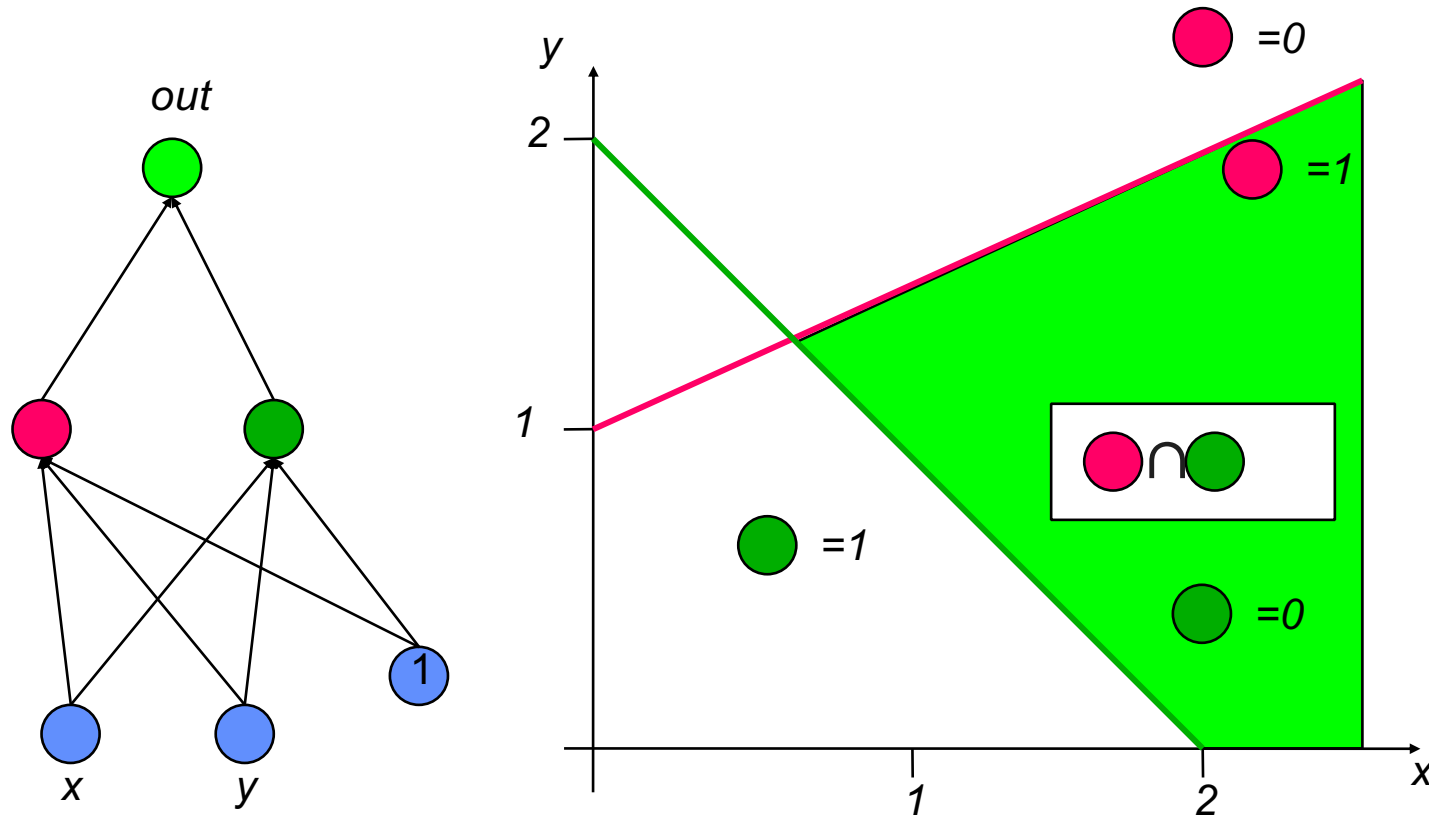
# Softmax

- Softmax layer as the output layer, to provide probabilistic interpretation.
- Output of network can be any value, which is difficult to interpret



# WHY Neural Network works?

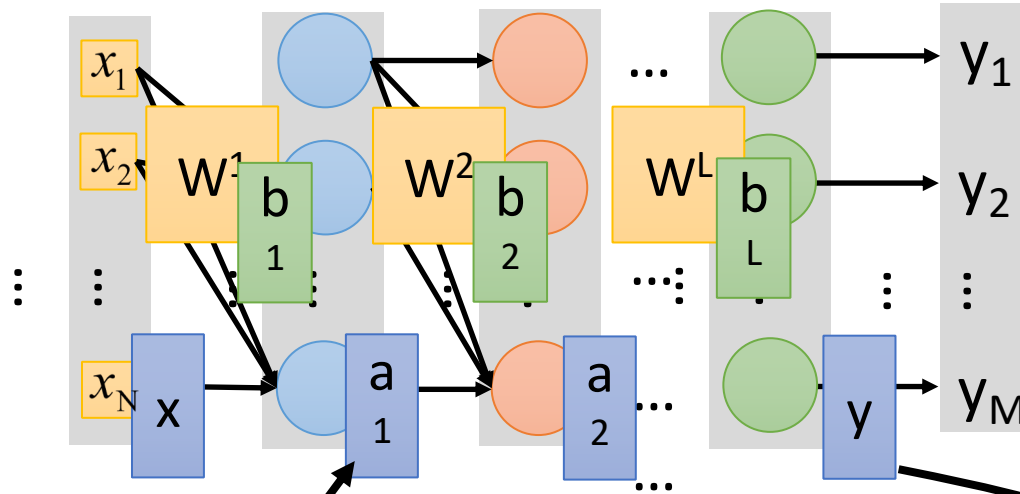
Visualizing as Constraint Satisfaction Networks without activation function



Each neuron in the hidden layer computes a nonlinear transform of inputs it receives. Hidden layer neurons act as “features” for final layer (a linear model) to produce output. The overall effect is a nonlinear mapping from inputs to outputs.

# WHY Neural Network works?

Why we require  
activation layers?  
Or  
Why we want to  
introduce non-  
linearity?



$$\sigma(W^1 x + b_1)$$

$$\sigma(W^L a_1^{L-1} + b_L)$$

$$y = f(x) = \sigma(W^L \dots \sigma(W^2 \sigma(W^1 x + b_1) + b_2) \dots + b_L)$$

It is related to kernel learning where kernels are learnt implicitly.  
What happen if we use linear mapping as an activation function?

MLP or NN with a single, sufficiently wide hidden layer can approximate any function.

## Preliminaries

**A visual proof that  
neural nets can  
approximate any  
function**

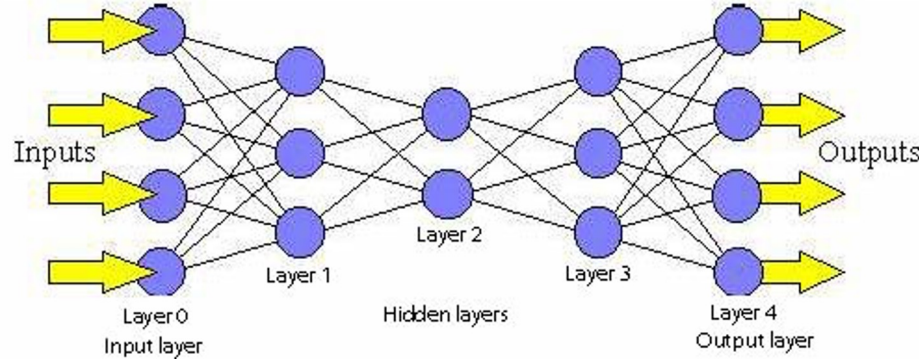
[http://neuralnetworksa  
nddeeplearning.com/c  
hap4.html](http://neuralnetworksanddeeplearning.com/chap4.html)

Why non-linearity is  
required?

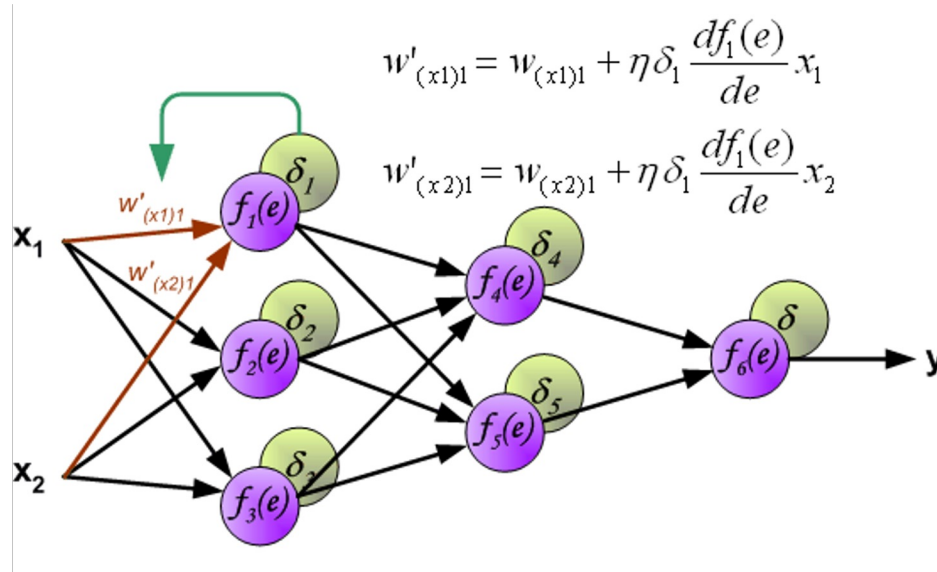
- 1) [Example of fitting non-linear data without non-linearities](#)
- 2) [Example of fitting non-linear data with non-linearities](#)



# Feed forward/Backpropagation Neural Network



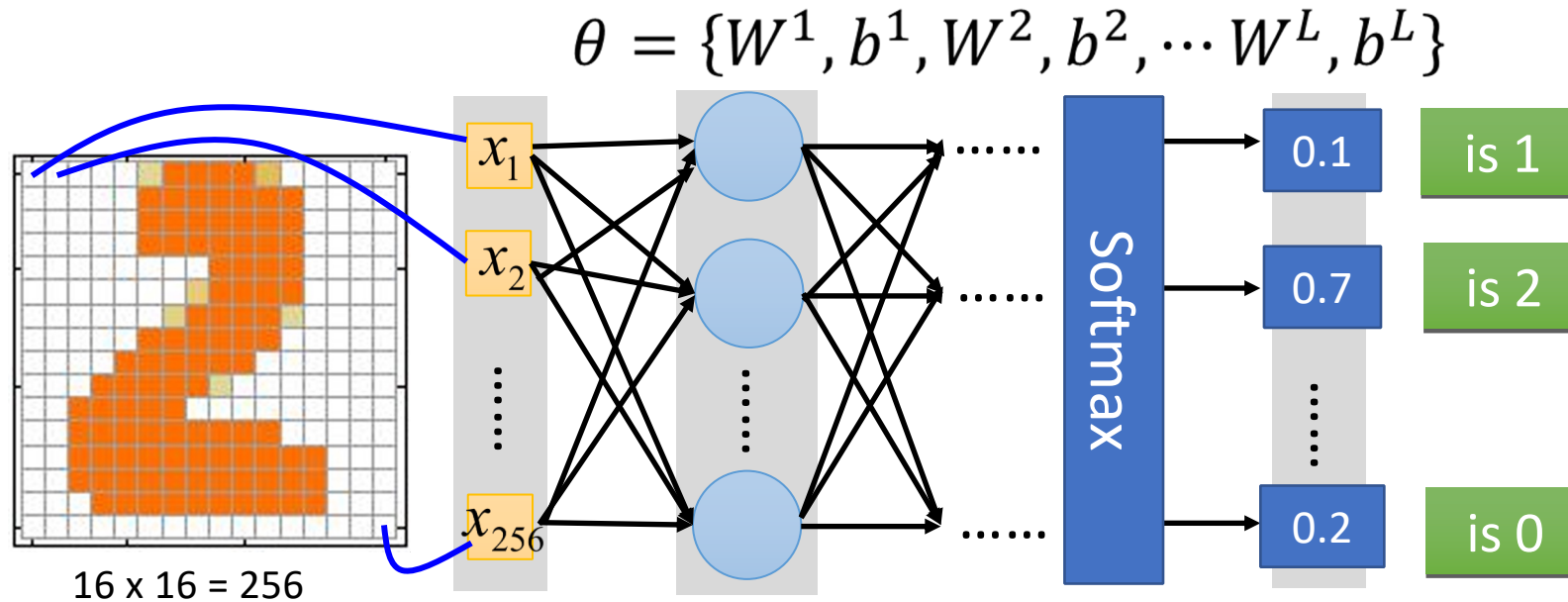
Feed forward algorithm



Backpropagation:



- Initialize the parameters
- Compute total error
- Then calculate gradient w.r.t each weight and eventually modify the weights to get better results.

# Learning network parameters



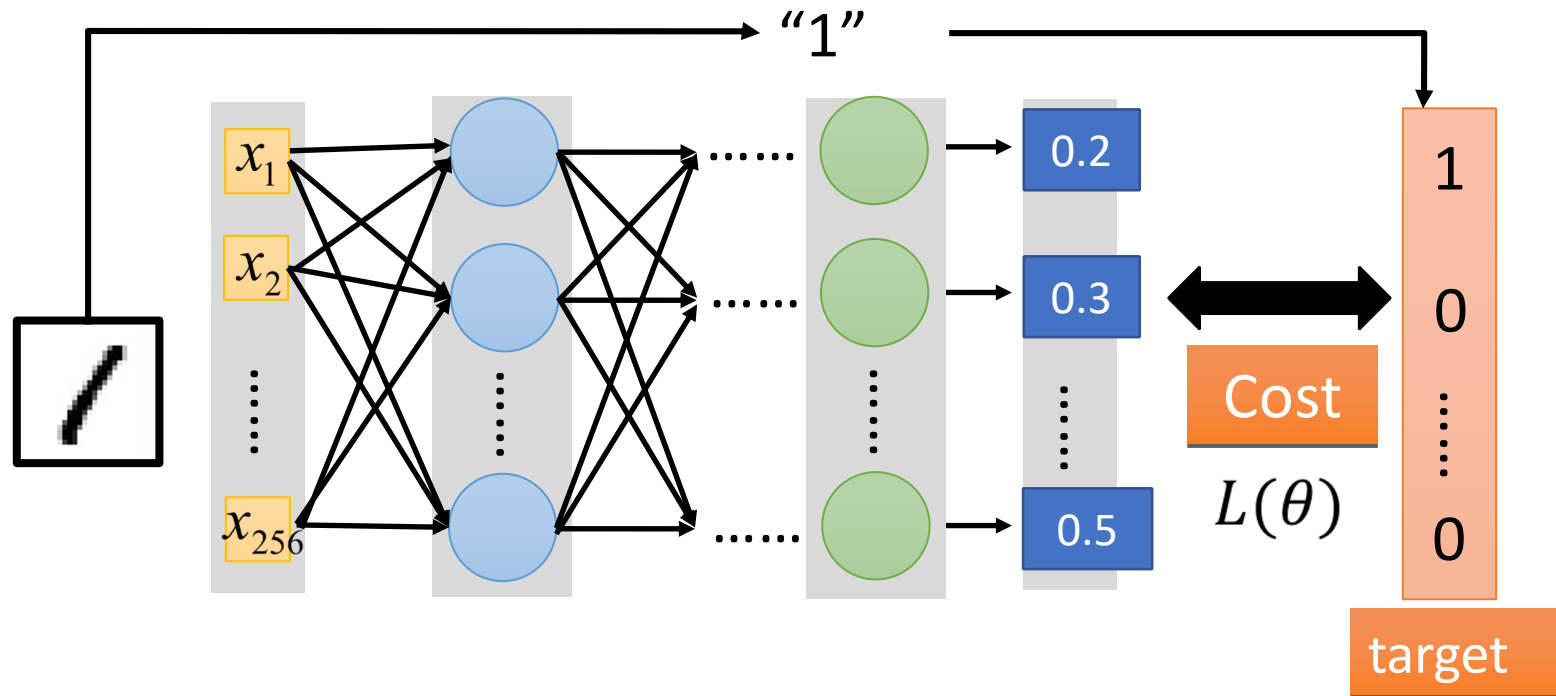
Orange  $\rightarrow$  1  
Black  $\rightarrow$  0

Aim: Find network parameters such that .....

Input   $\rightarrow$   $y_1$  has the maximum value  
:  
Input   $\rightarrow$   $y_2$  has the maximum value  
:

How can the neural network achieve this ?

# Computing cost of a sample



Cost can be Euclidean distance or cross entropy between network output and target  
Given a set of network parameters, each example has a cost value.

How to evaluate that  
how bad the network  
parameters performs on  
this task?

Total Cost: 
$$C(\theta) = \sum_{r=1}^R L^r(\theta)$$

Find the network parameters that minimize total cost

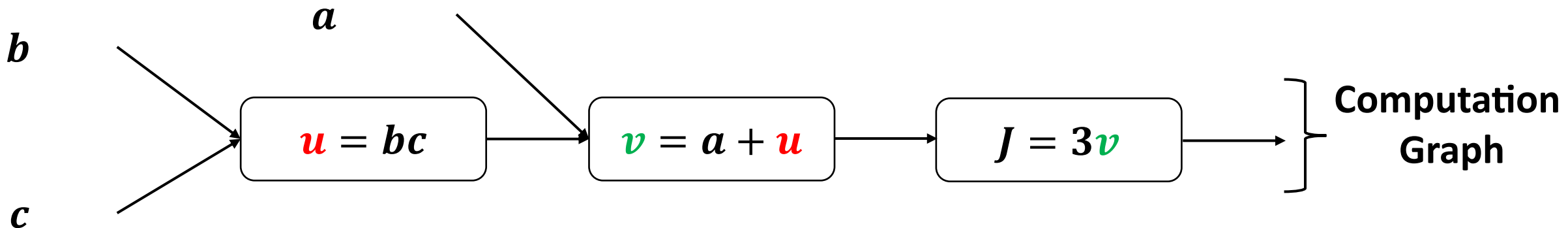
# The Flow of Computations in Neural Networks

- 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”
  - What is a computation graph?

# What is a Computation Graph?

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

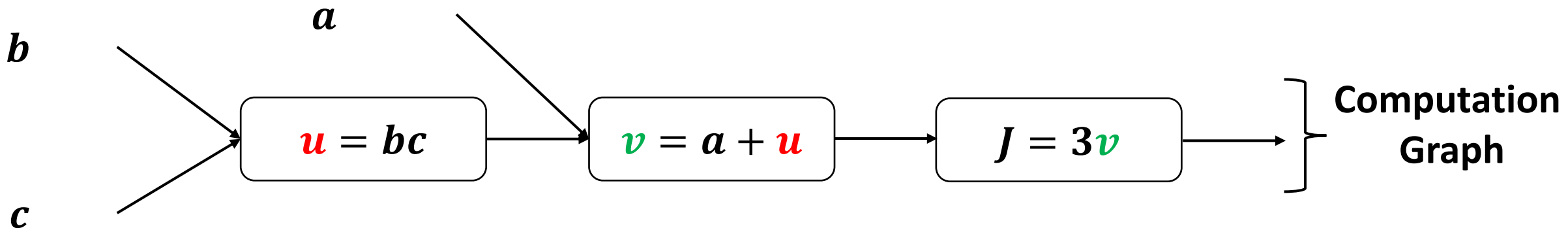
$$\underbrace{\underbrace{\underbrace{J(a, b, c) = 3(a + \underbrace{bc}_{\textcolor{red}{u}})}_{\textcolor{green}{v}}}_{J}} \quad \left. \vphantom{\underbrace{\underbrace{\underbrace{J(a, b, c) = 3(a + \underbrace{bc}_{\textcolor{red}{u}})}_{\textcolor{green}{v}}}_{J}}} \right\} \begin{array}{l} \textcolor{red}{u} = bc \\ \textcolor{green}{v} = a + \textcolor{red}{u} \\ J = 3\textcolor{green}{v} \end{array}$$



# Forward Propagation

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

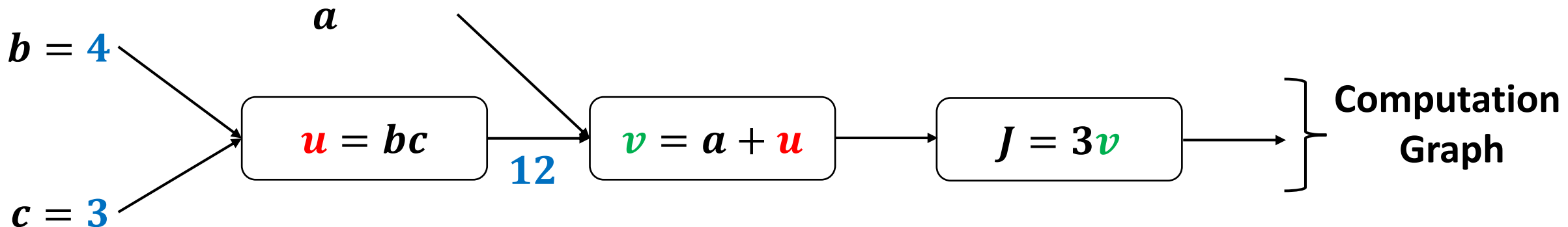
$$\underbrace{\underbrace{\underbrace{J(a, b, c) = 3(a + \underbrace{bc}_{\textcolor{red}{u}})}_{\textcolor{green}{v}}}_{J}} \quad \left. \vphantom{\underbrace{\underbrace{\underbrace{J(a, b, c) = 3(a + \underbrace{bc}_{\textcolor{red}{u}})}_{\textcolor{green}{v}}}_{J}}} \right\} \begin{array}{l} \textcolor{red}{u} = bc \\ \textcolor{green}{v} = a + \textcolor{red}{u} \\ J = 3\textcolor{green}{v} \end{array}$$



# Forward Propagation

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

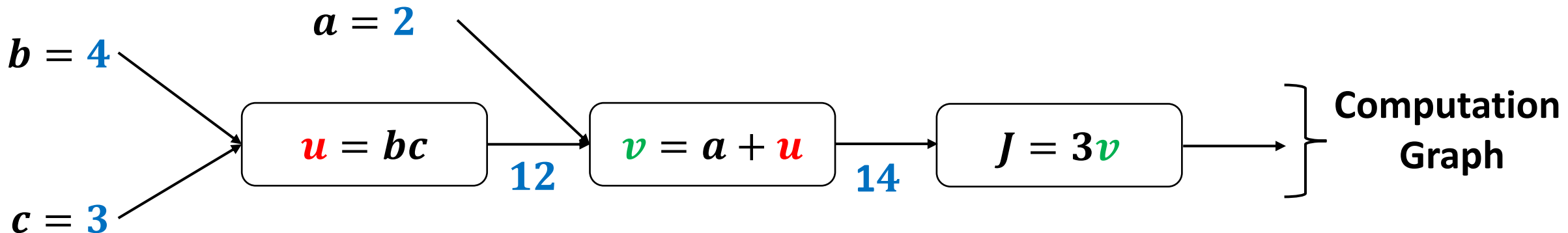
$$\underbrace{\underbrace{\underbrace{J(a, b, c) = 3(a + \underbrace{bc}_{\textcolor{red}{u}})}_{\textcolor{green}{v}}}_{J}} \quad \left. \vphantom{\underbrace{\underbrace{\underbrace{J(a, b, c) = 3(a + \underbrace{bc}_{\textcolor{red}{u}})}_{\textcolor{green}{v}}}_{J}}} \right\} \begin{array}{l} \textcolor{red}{u} = bc \\ \textcolor{green}{v} = a + \textcolor{red}{u} \\ J = 3\textcolor{green}{v} \end{array}$$



# Forward Propagation

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

$$\underbrace{\underbrace{\underbrace{J(a, b, c) = 3(a + \underbrace{bc}_{\textcolor{red}{u}})}_{\textcolor{green}{v}}}_{J}} \quad \left. \vphantom{\underbrace{\underbrace{\underbrace{J(a, b, c) = 3(a + \underbrace{bc}_{\textcolor{red}{u}})}_{\textcolor{green}{v}}}_{J}}} \right\} \begin{array}{l} \textcolor{red}{u} = bc \\ \textcolor{green}{v} = a + \textcolor{red}{u} \\ J = 3\textcolor{green}{v} \end{array}$$

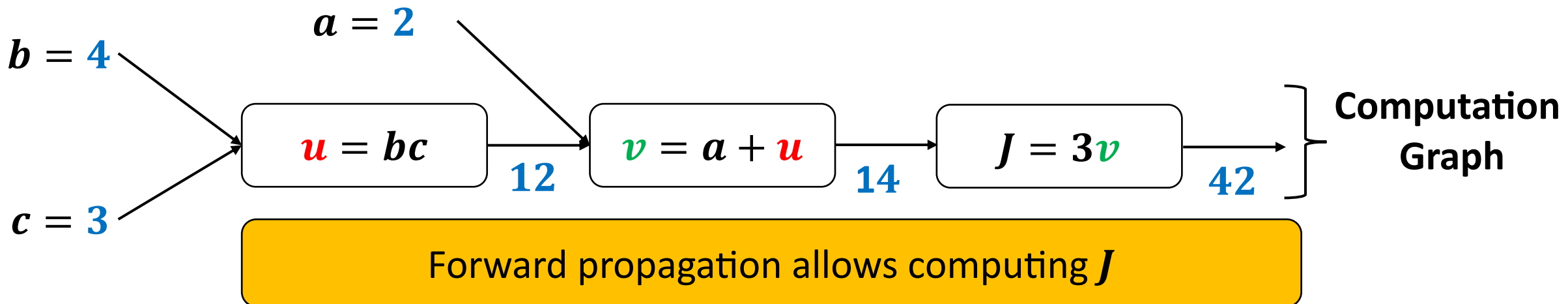




# Forward Propagation

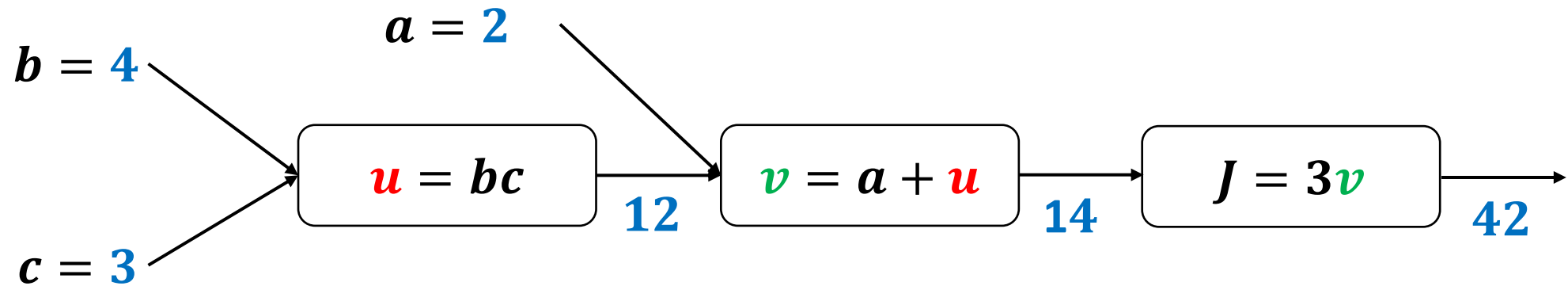
- Let us assume we want to compute the following function  $J$ :

$$\begin{array}{c} J(a, b, c) = 3(a + \underbrace{bc}_{\underbrace{u}_{v}}) \\ \underbrace{\hspace{1.5cm}}_J \end{array} \left. \begin{array}{l} \textcolor{red}{u} = bc \\ \textcolor{green}{v} = a + \textcolor{red}{u} \\ J = 3\textcolor{green}{v} \end{array} \right\}$$



# Backward Propagation

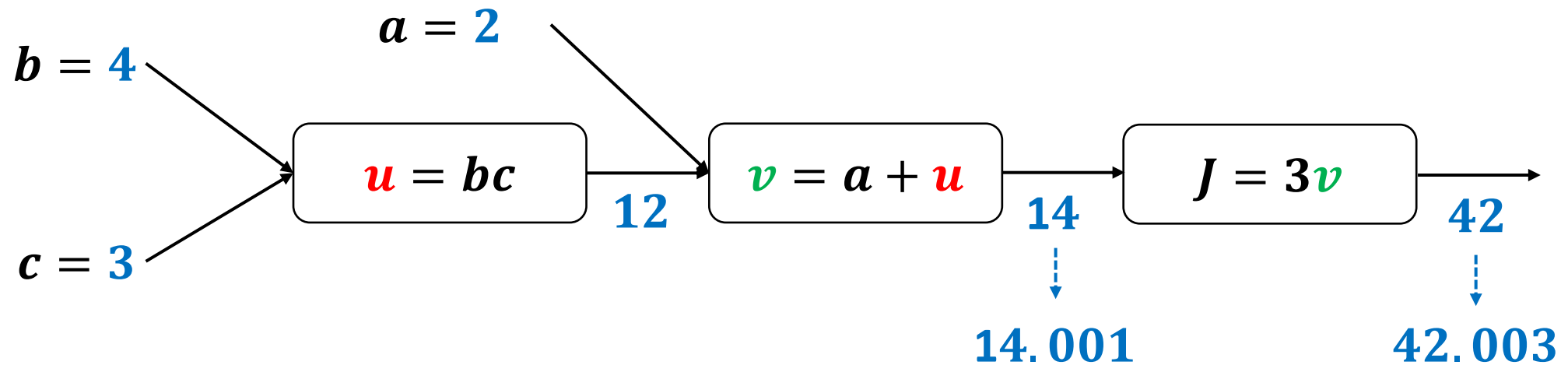
- Let us now compute the derivatives of the variables through the computation graph as follows:



$$\frac{dJ}{dv} = \text{Derivative of } J \text{ with respect to } v$$

# Backward Propagation

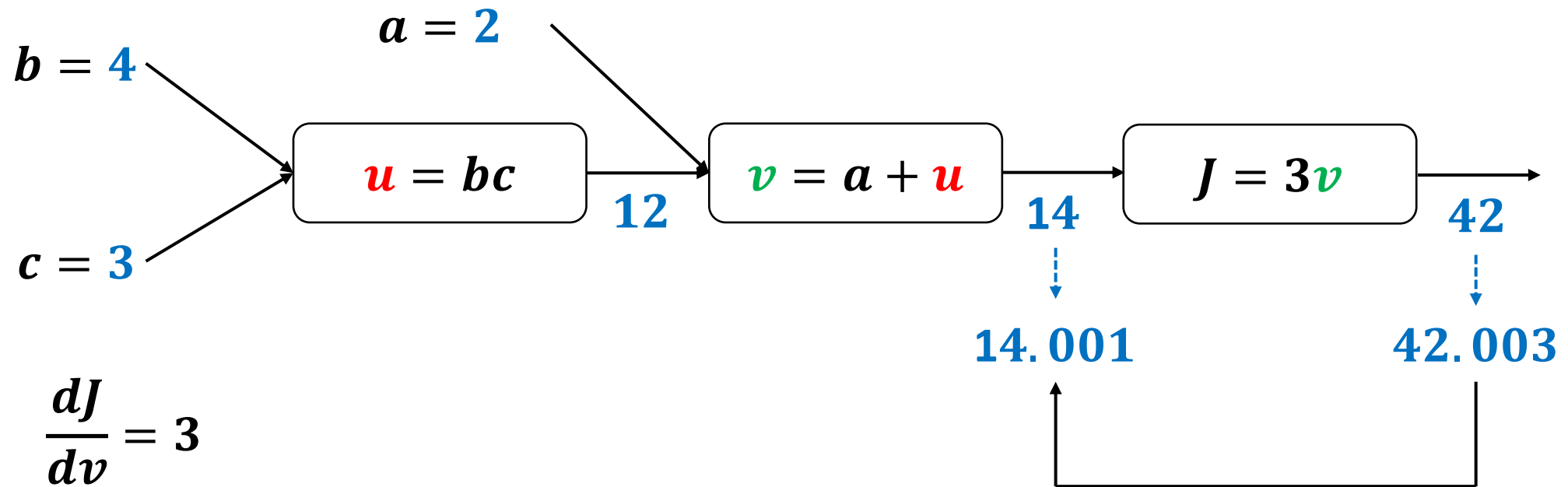
- Let us now compute the derivatives of the variables through the computation graph as follows:



$\frac{dJ}{dv}$  = If we change  $v$  a little bit, how would  $J$  change?

# Backward Propagation

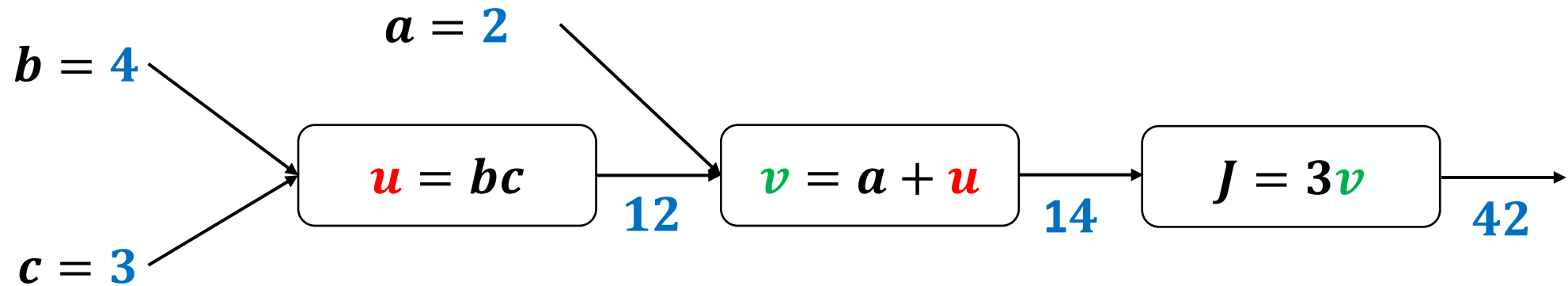
- Let us now compute the derivatives of the variables through the computation graph as follows:



To compute the derivative of  $J$  with respect to  $v$ , we went *back* to  $v$ , nudged it, and measured the corresponding resultant increase on  $J$

# Backward Propagation

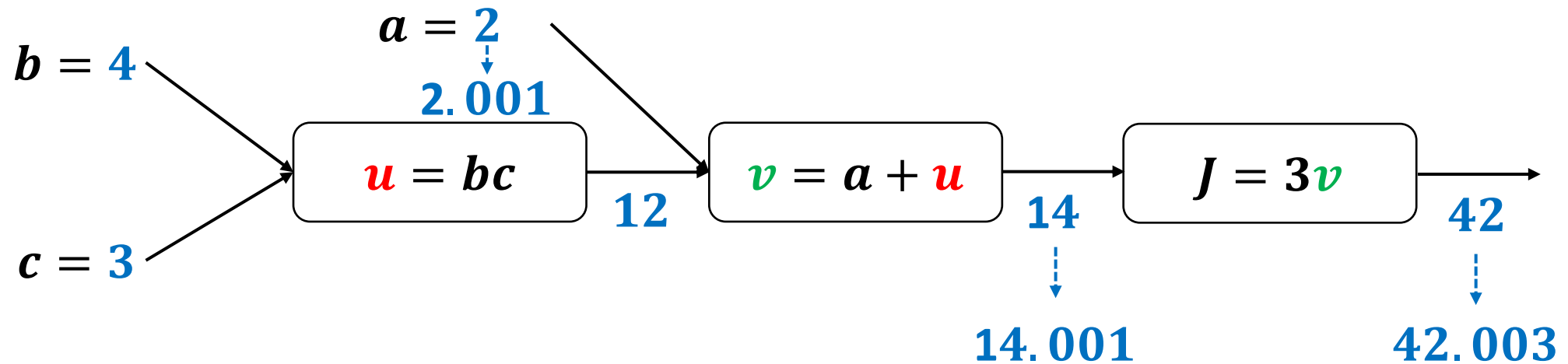
- Let us now compute the derivatives of the variables through the computation graph as follows:



$$\frac{dJ}{da} = \text{Derivative of } J \text{ with respect to } a$$

# Backward Propagation

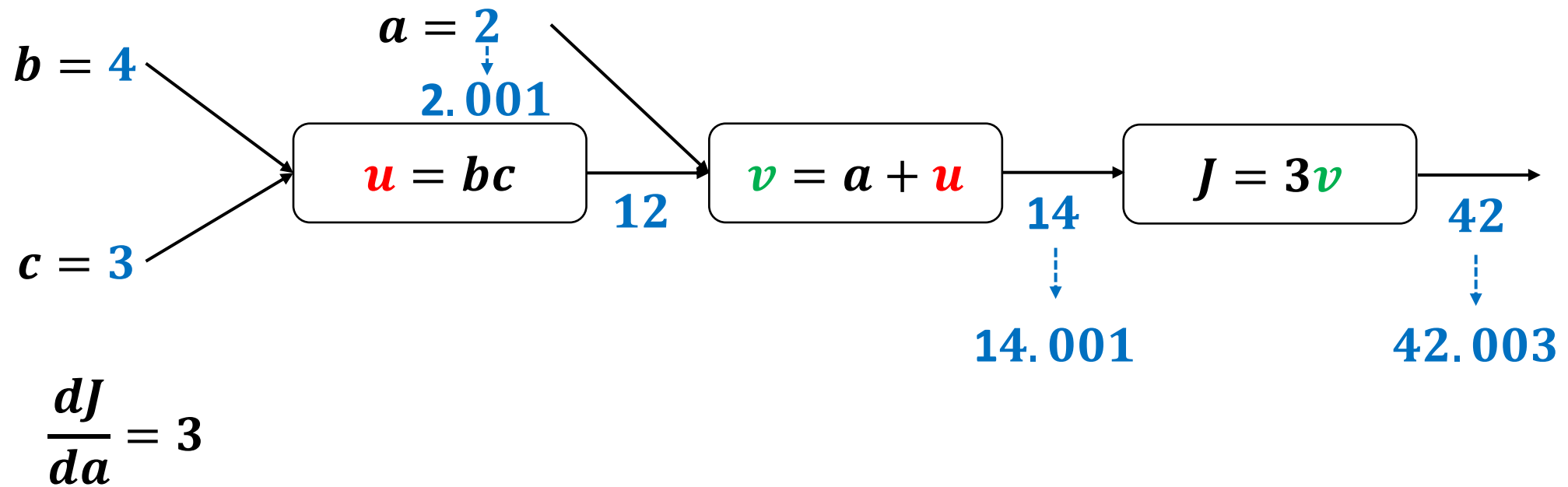
- Let us now compute the derivatives of the variables through the computation graph as follows:



$\frac{dJ}{da}$  = If we change  $a$  a little bit, how would  $J$  change?

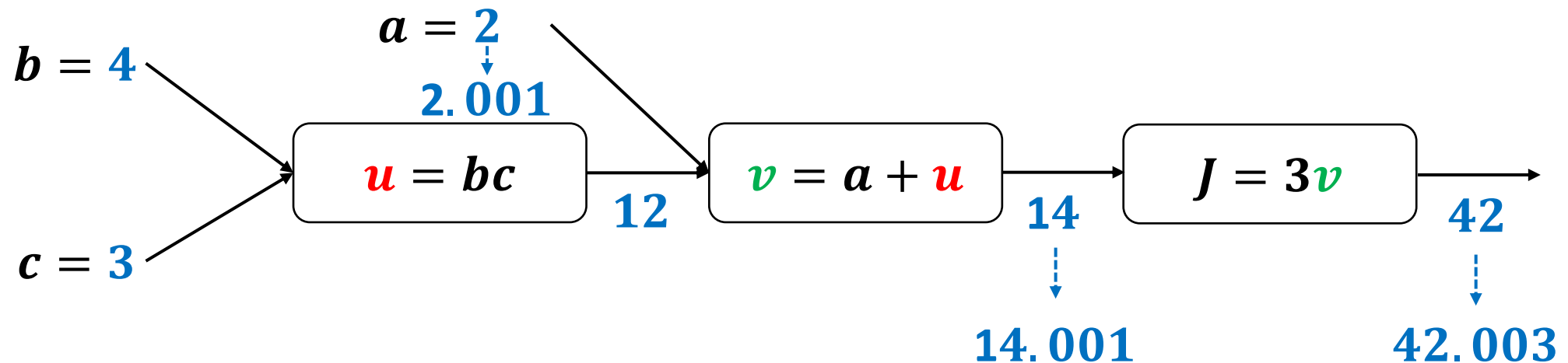
# Backward Propagation

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



# Backward Propagation

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



$$\frac{dJ}{da} = 3 =$$

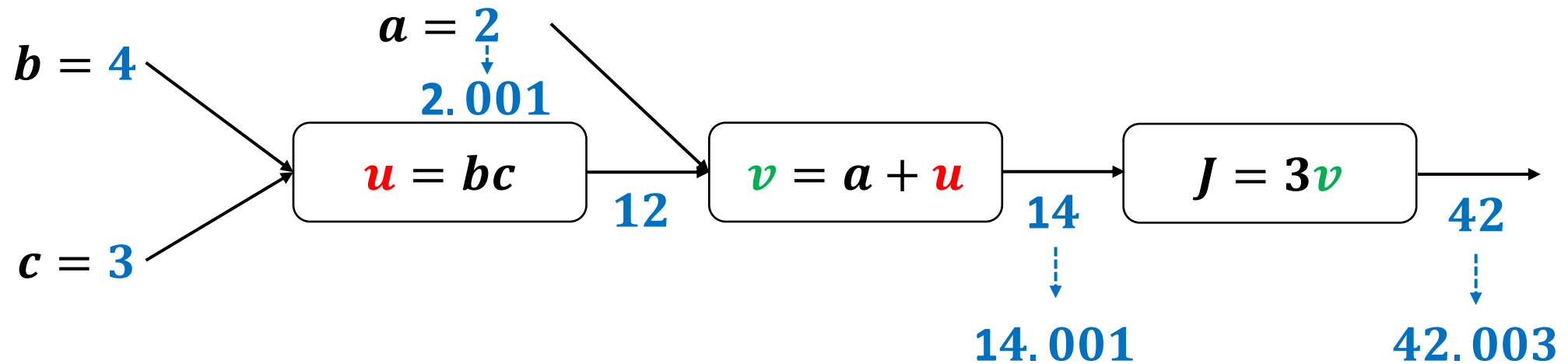
$$\frac{dv}{da}$$

The change in  $a$  caused a change in  $v$



# Backward Propagation

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

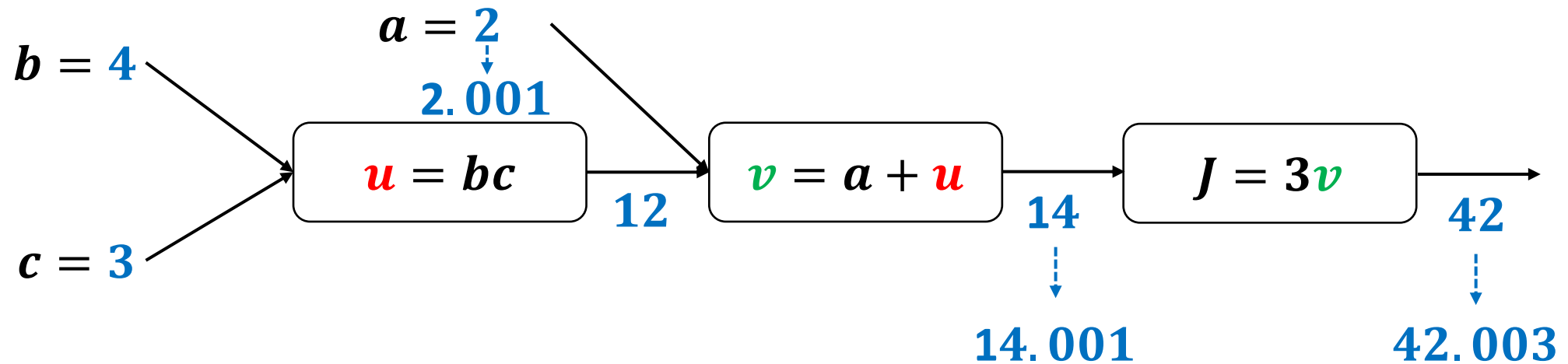


$$\frac{dJ}{da} = 3 = \frac{dJ}{dv} \times \frac{dv}{da}$$

And the change in  $v$  caused a change in  $J$

# Backward Propagation

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

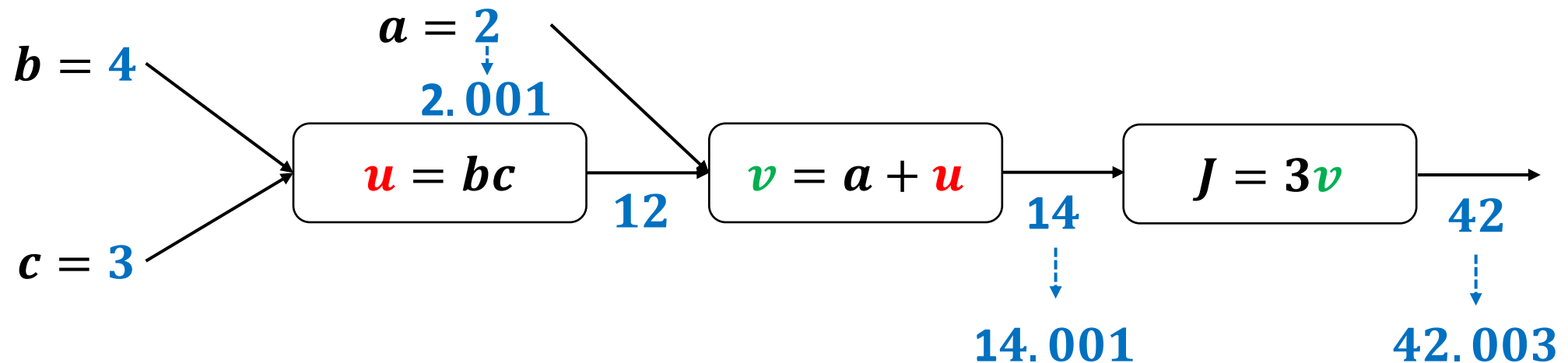


$$\frac{dJ}{da} = 3 = \frac{dJ}{dv} \times \frac{dv}{da}$$

This is denoted as *the chain rule* in calculus

# Backward Propagation

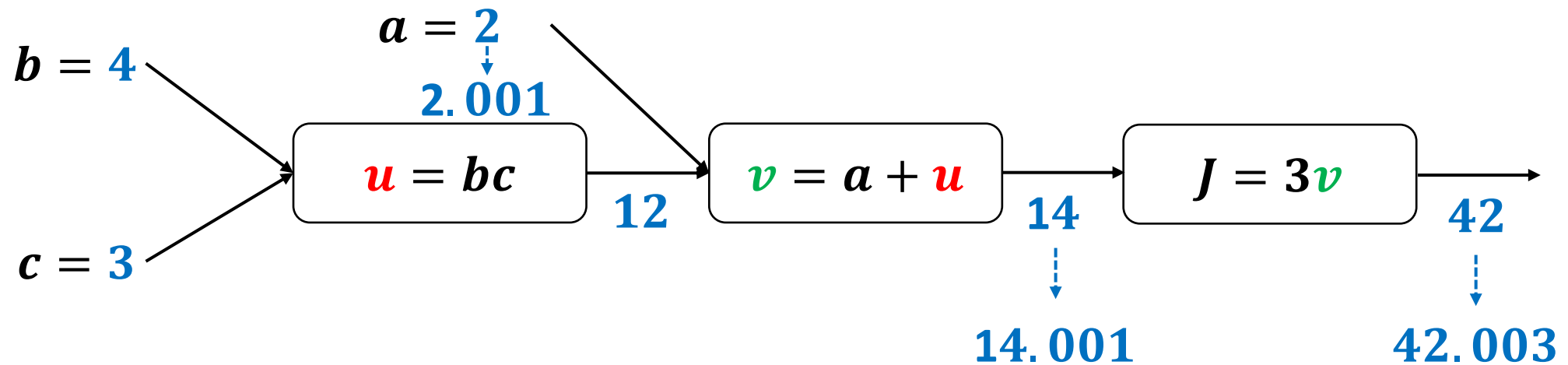
- Let us now compute the derivatives of the variables through the computation graph as follows:



$$\frac{dJ}{da} = 3 = \frac{dJ}{dv} \times 1$$

# Backward Propagation

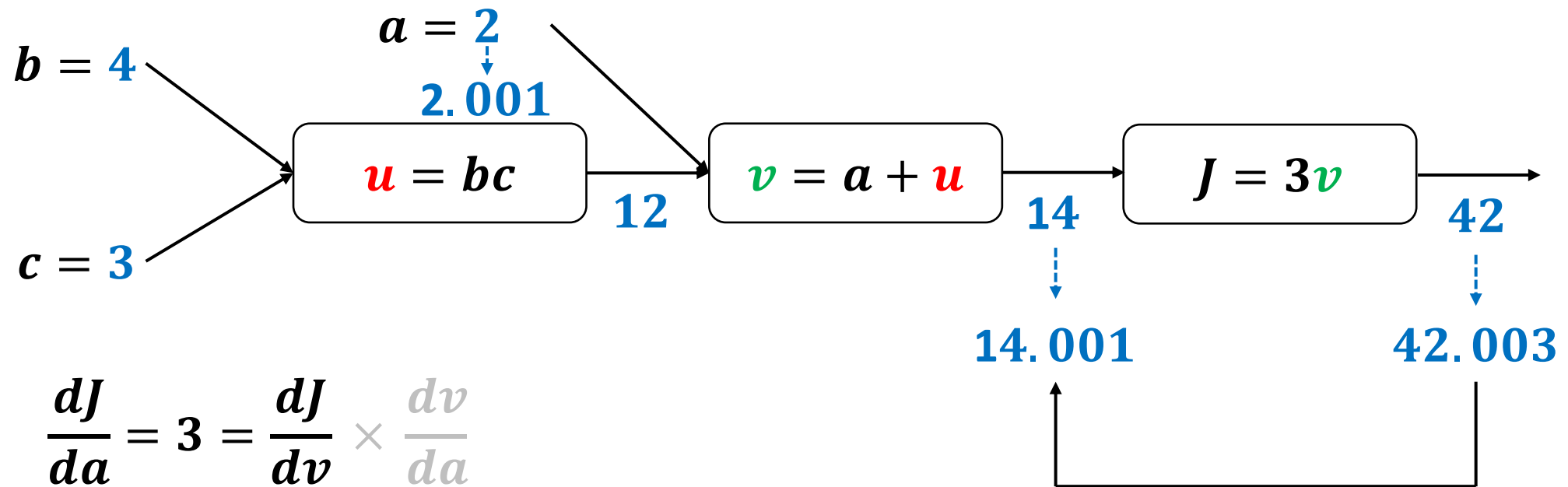
- Let us now compute the derivatives of the variables through the computation graph as follows:



$$\frac{dJ}{da} = 3 = 3 \times 1$$

# Backward Propagation

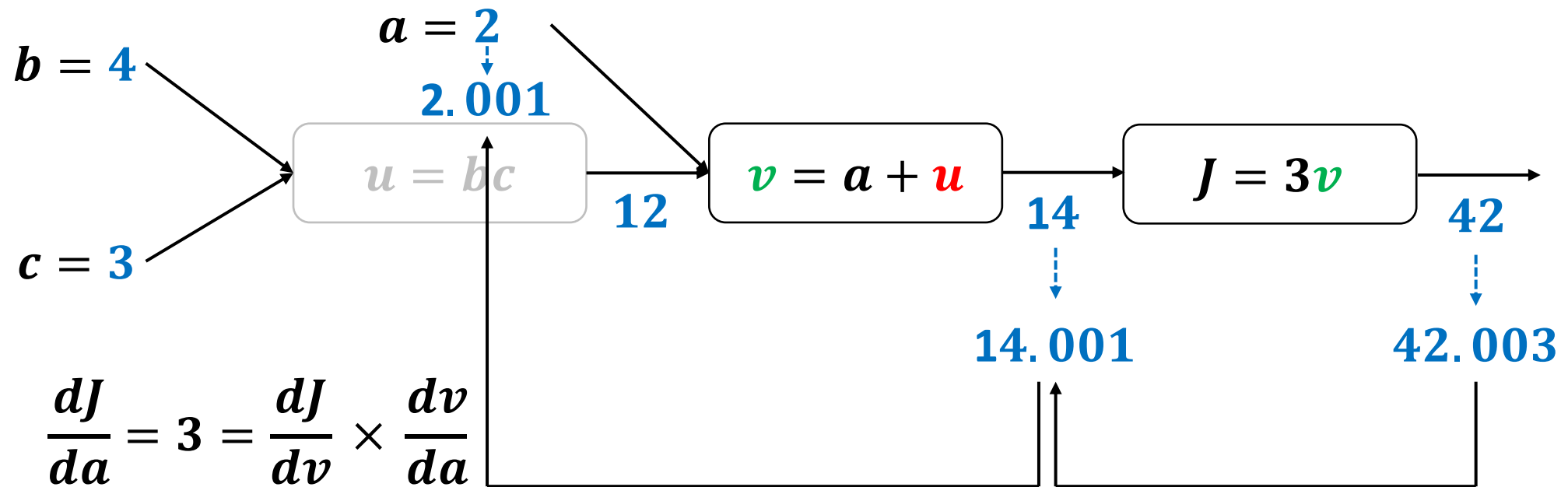
- Let us now compute the derivatives of the variables through the computation graph as follows:



In essence, to compute the derivative of  $J$  with respect to  $a$ , we had to go *back* to  $v$ , nudge it a little bit, and measure the corresponding resultant increase on  $J$

# Backward Propagation

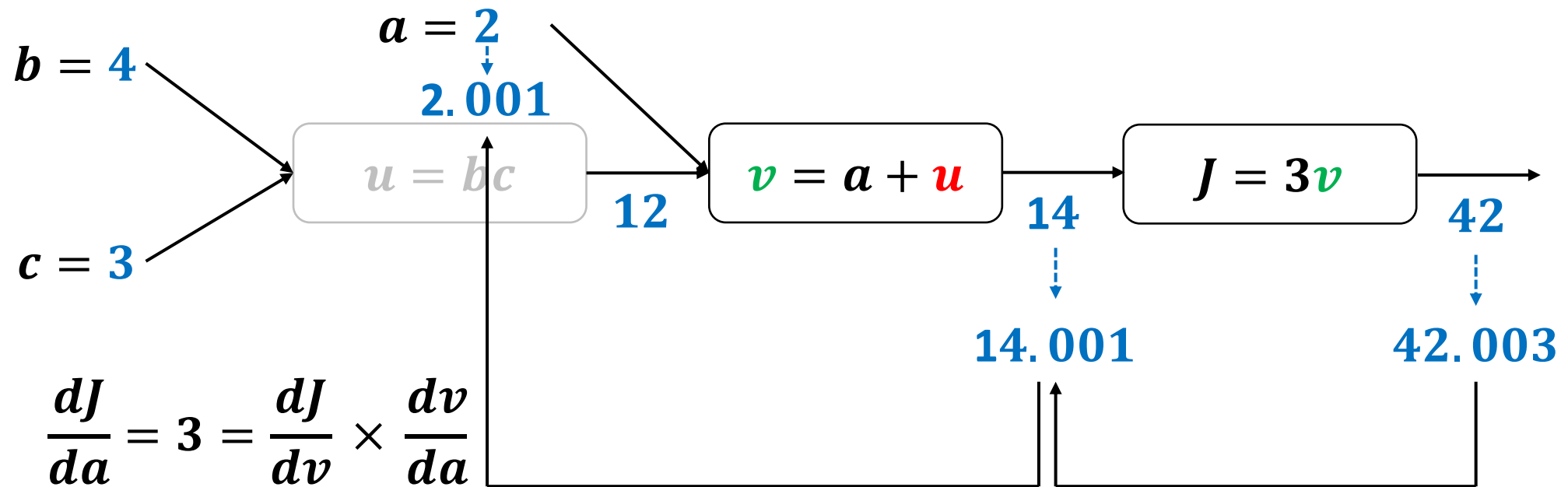
- Let us now compute the derivatives of the variables through the computation graph as follows:



Then, we had to go *back* to  $a$ , nudge it a little bit, and measure the corresponding resultant increase on  $v$

# Backward Propagation

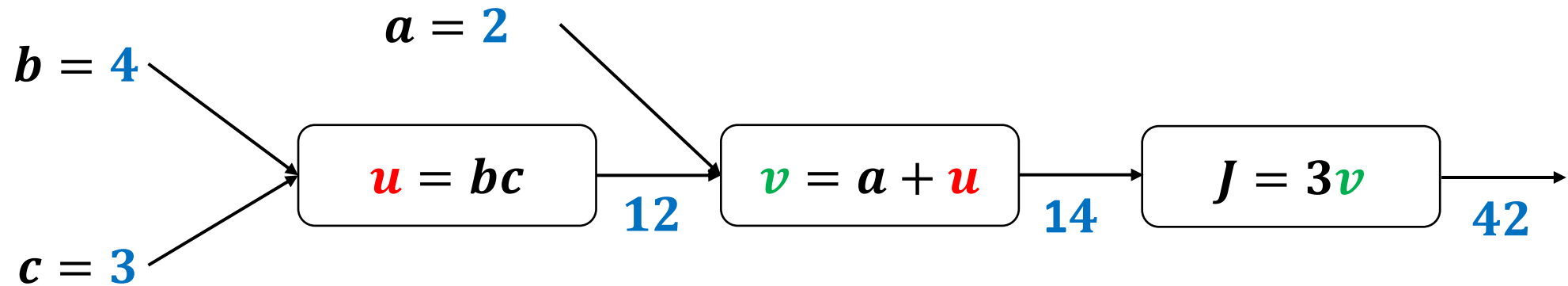
- Let us now compute the derivatives of the variables through the computation graph as follows:



Then, we multiplied the changes together (i.e., we applied the chain rule!)

# Backward Propagation

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

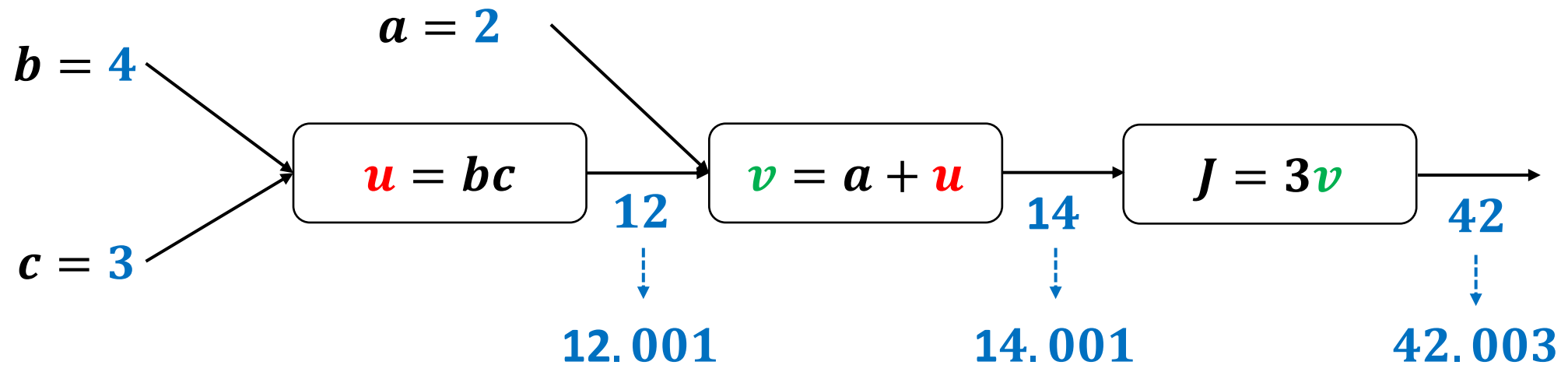


$\frac{dJ}{du}$  = Derivative of  $J$  with respect to  $u$



# Backward Propagation

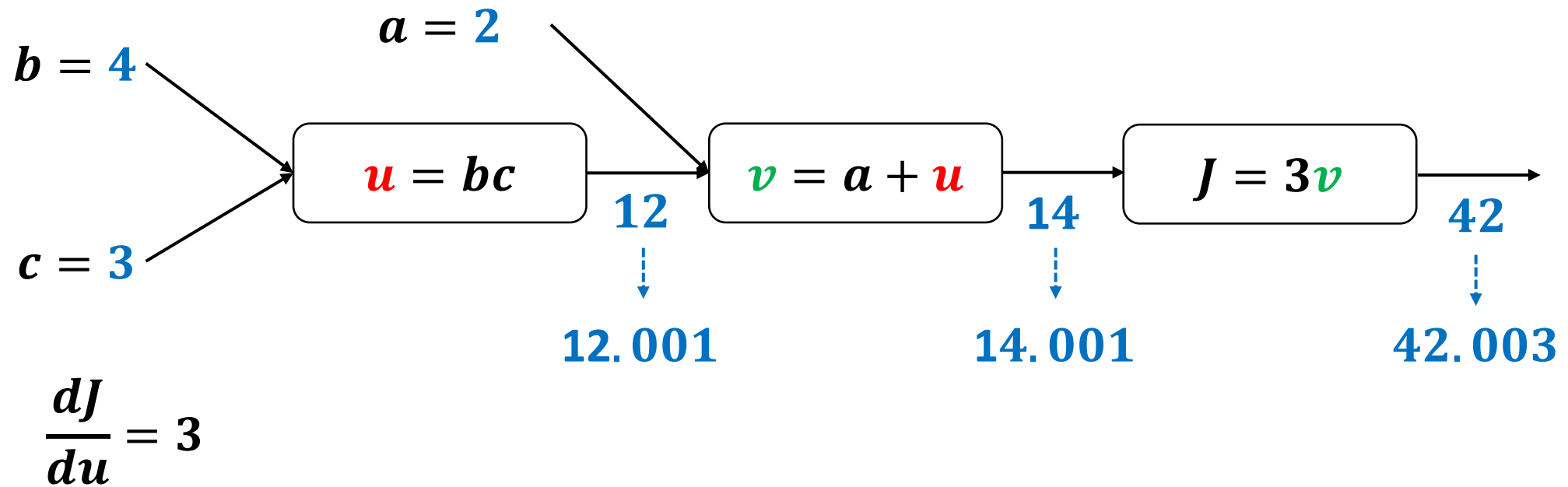
- Let us now compute the derivatives of the variables through the computation graph as follows:



$\frac{dJ}{du}$  = If we change  $u$  a little bit, how would  $J$  change?

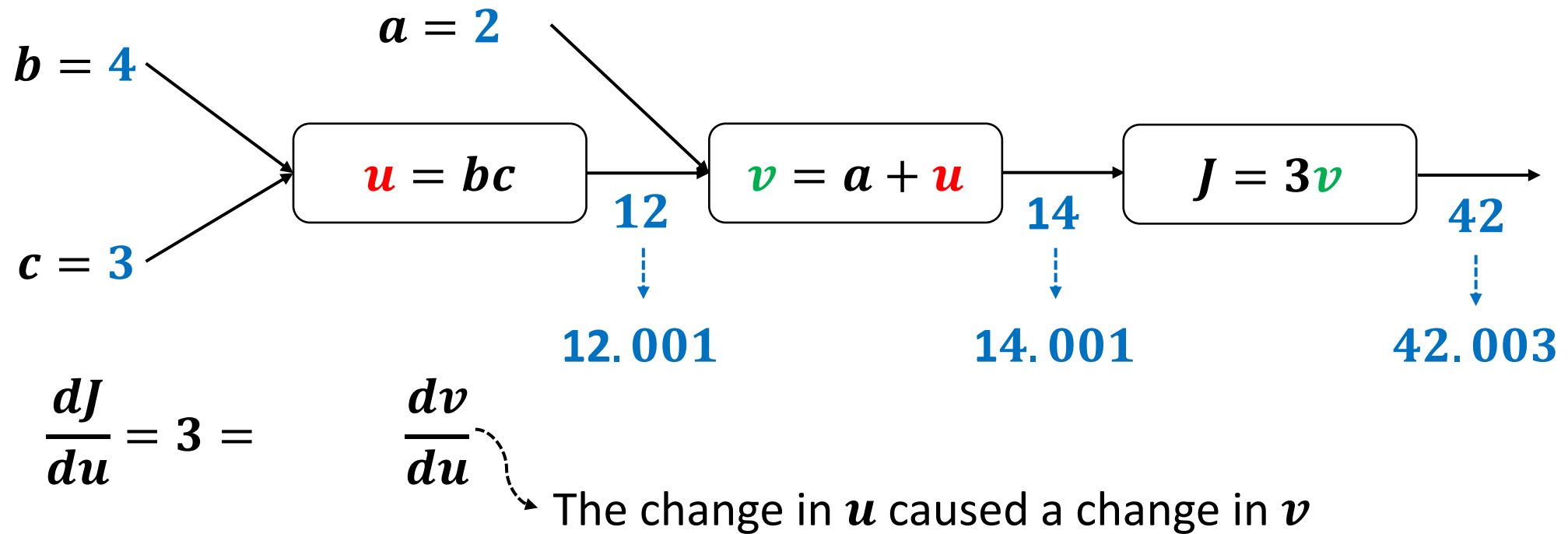
# Backward Propagation

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



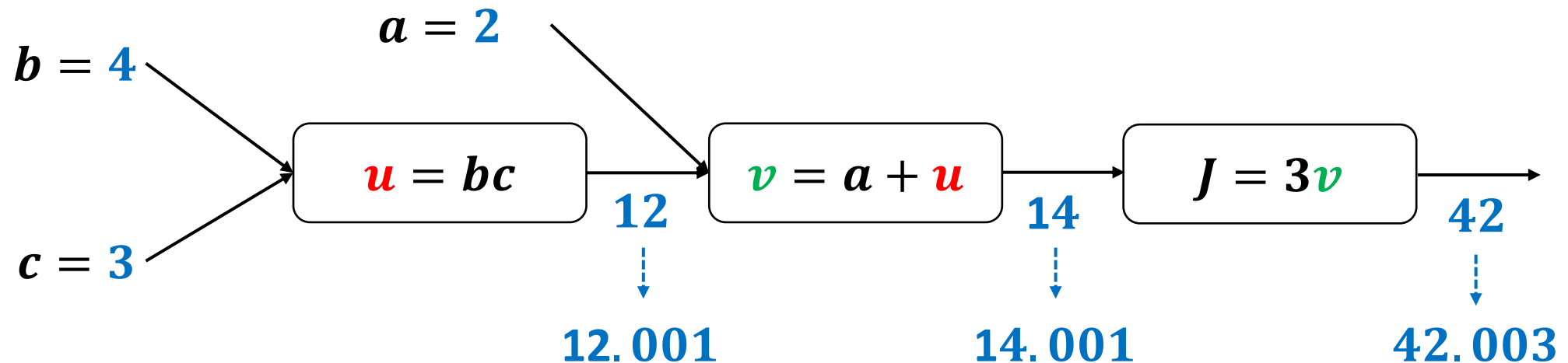
# Backward Propagation

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



# Backward Propagation

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

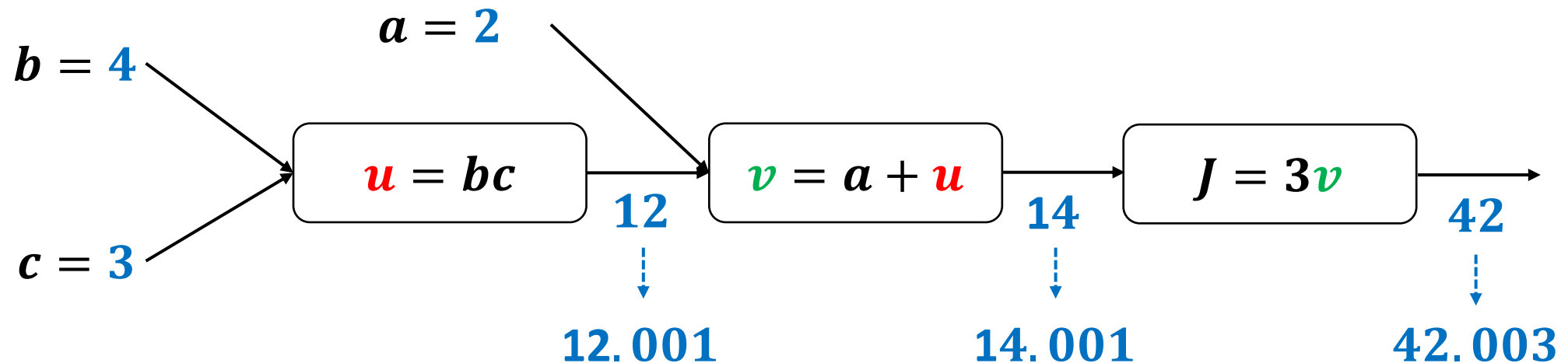


$$\frac{dJ}{du} = 3 = \frac{dJ}{dv} \times \frac{dv}{du}$$

And the change in  $v$  caused a change in  $J$

# Backward Propagation

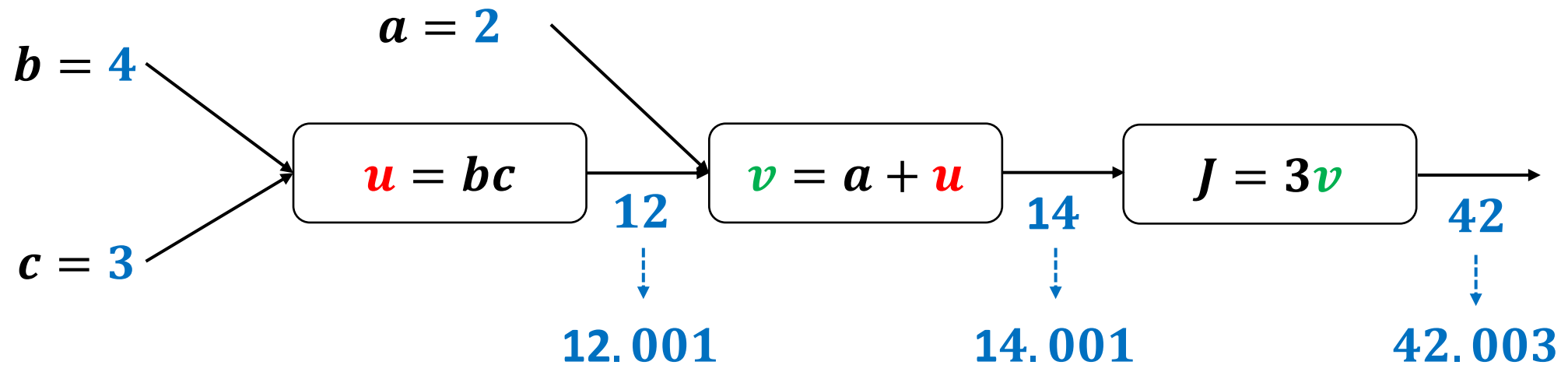
- Let us now compute the derivatives of the variables through the computation graph as follows:



$$\frac{dJ}{du} = 3 = \frac{dJ}{dv} \times 1$$

# Backward Propagation

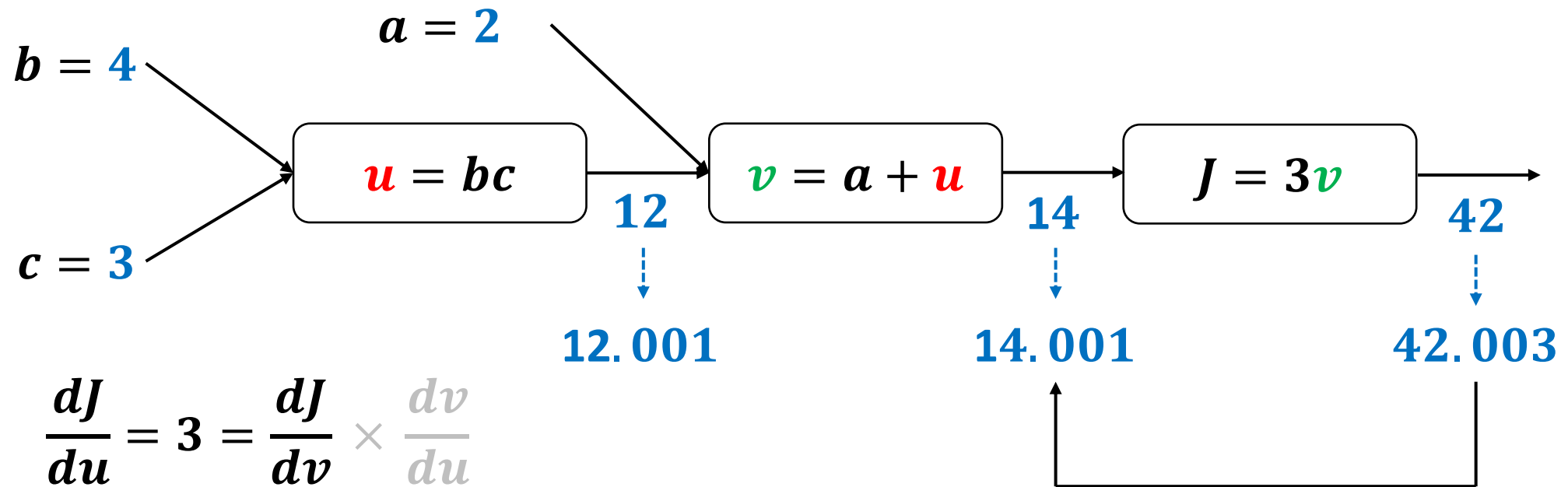
- Let us now compute the derivatives of the variables through the computation graph as follows:



$$\frac{dJ}{du} = 3 = 3 \times 1$$

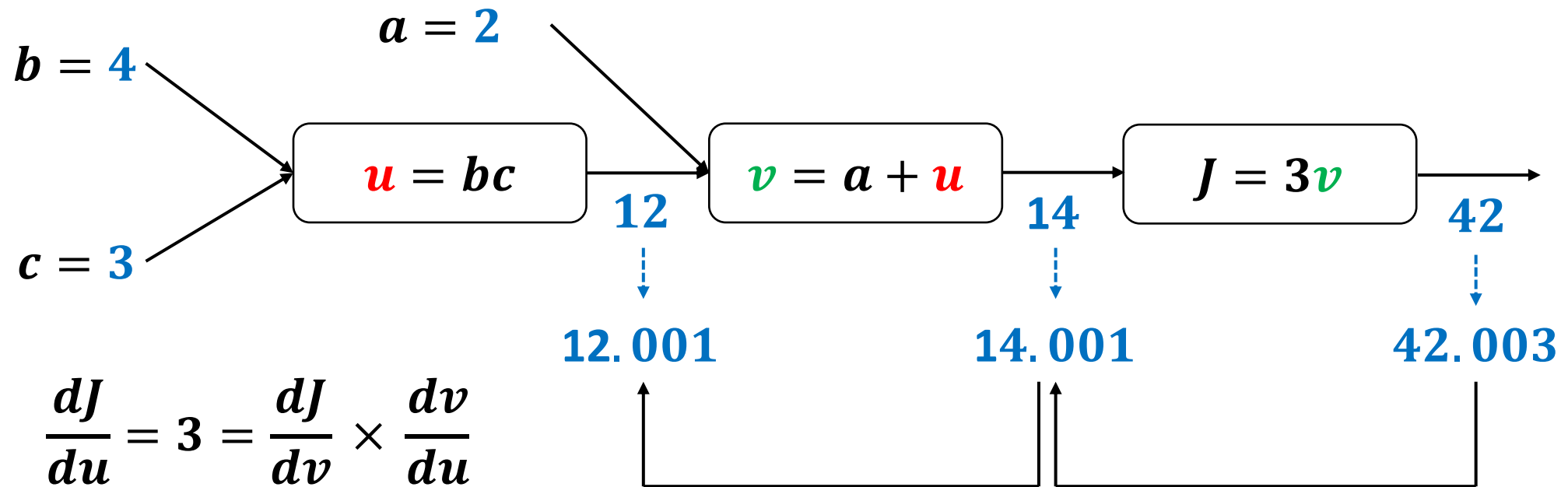
# Backward Propagation

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



# Backward Propagation

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

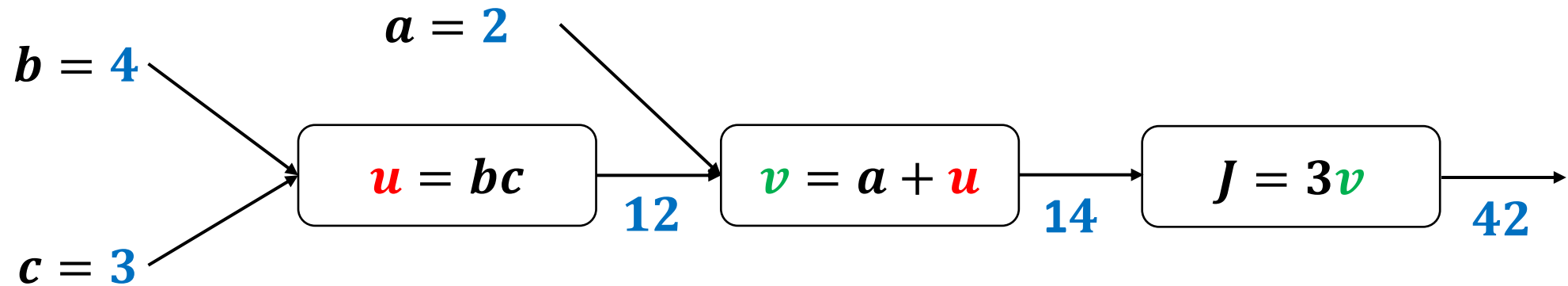


Same as before, we had to go back to  $v$  then to  $u$  in order to compute the derivative of  $J$



# Backward Propagation

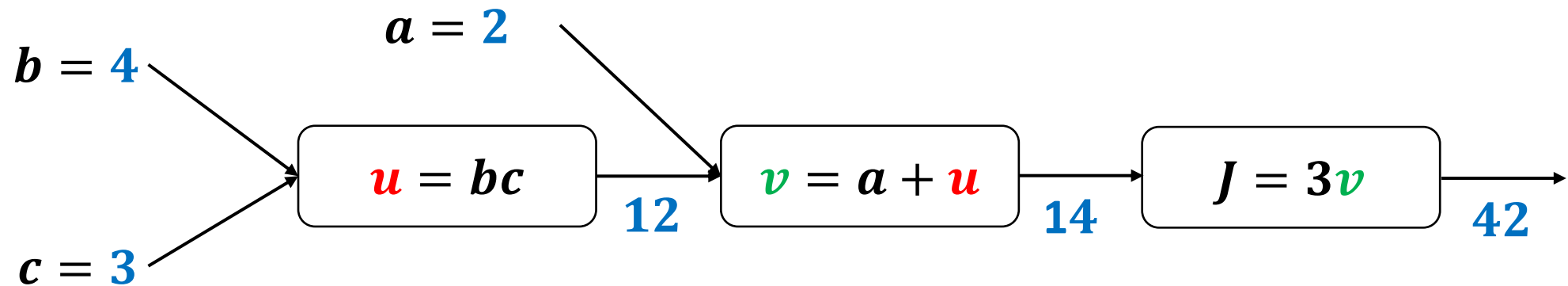
- Let us now compute the derivatives of the variables through the computation graph as follows:



$$\frac{dJ}{db} = \text{Derivative of } J \text{ with respect to } b$$

# Backward Propagation

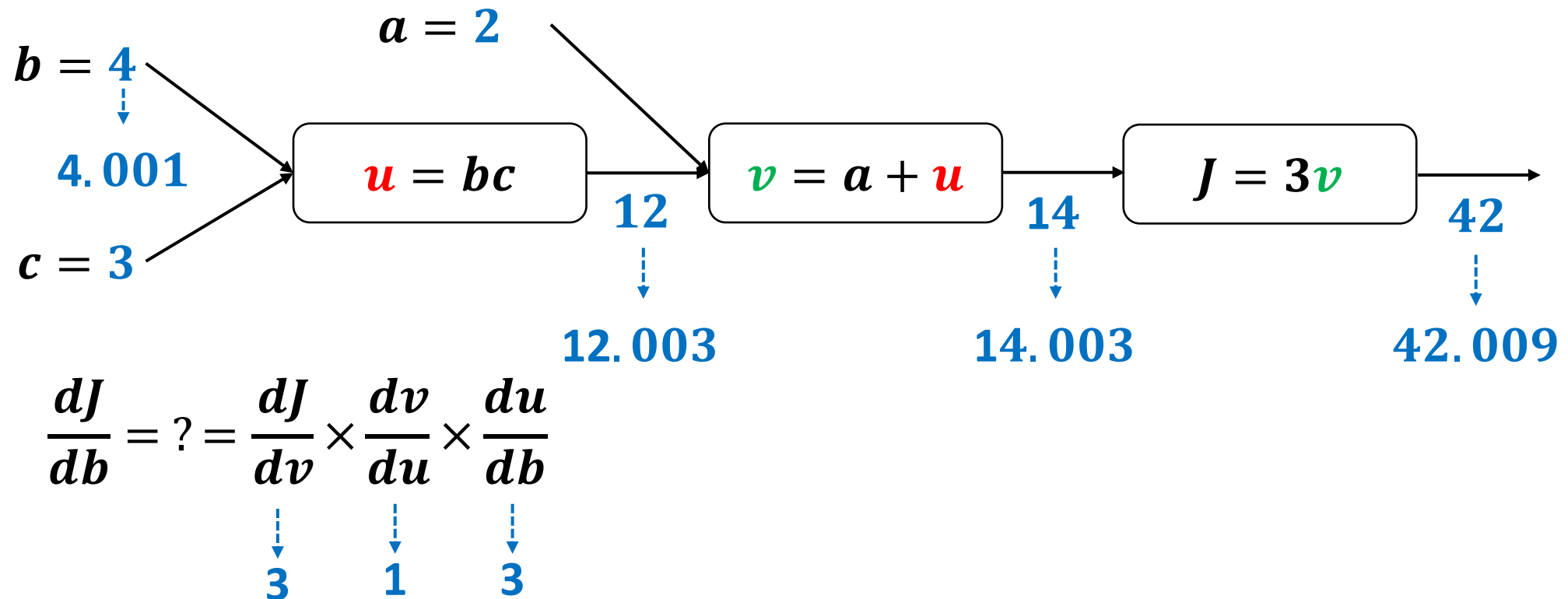
- Let us now compute the derivatives of the variables through the computation graph as follows:



$\frac{dJ}{db}$  = If we change  $b$  a little bit, how would  $J$  change?

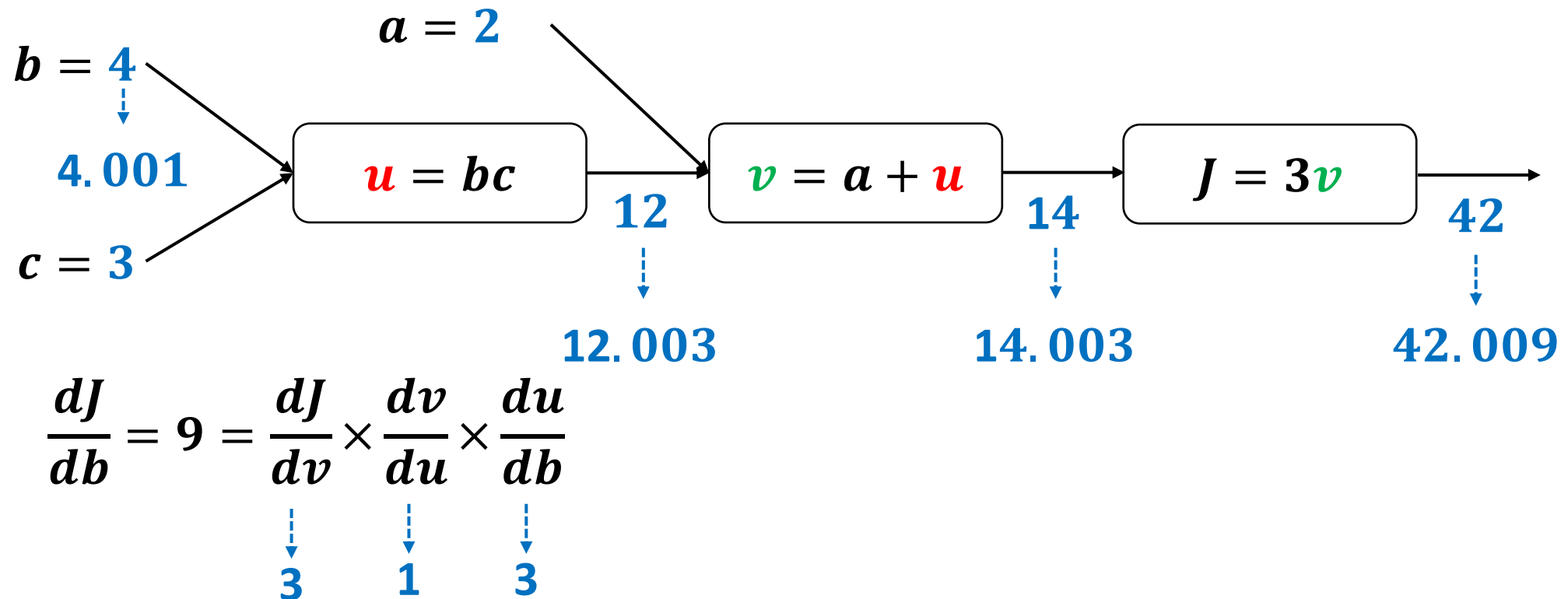
# Backward Propagation

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



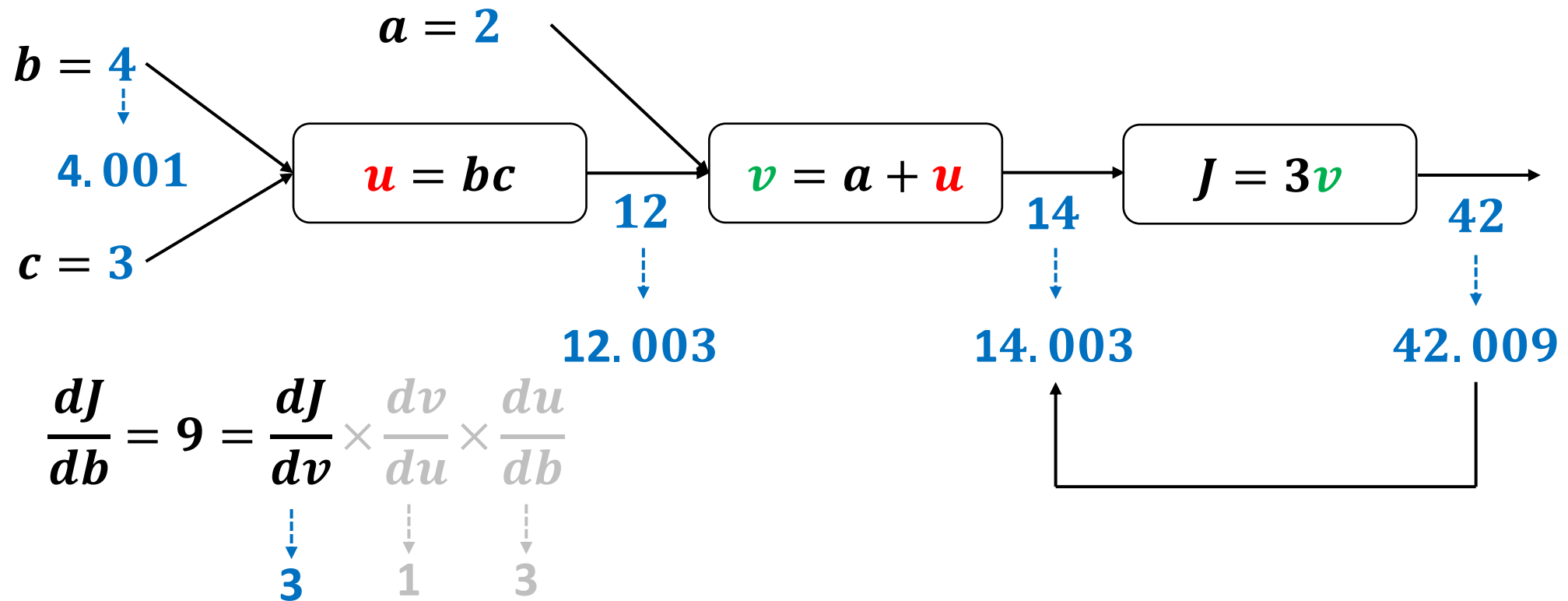
# Backward Propagation

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



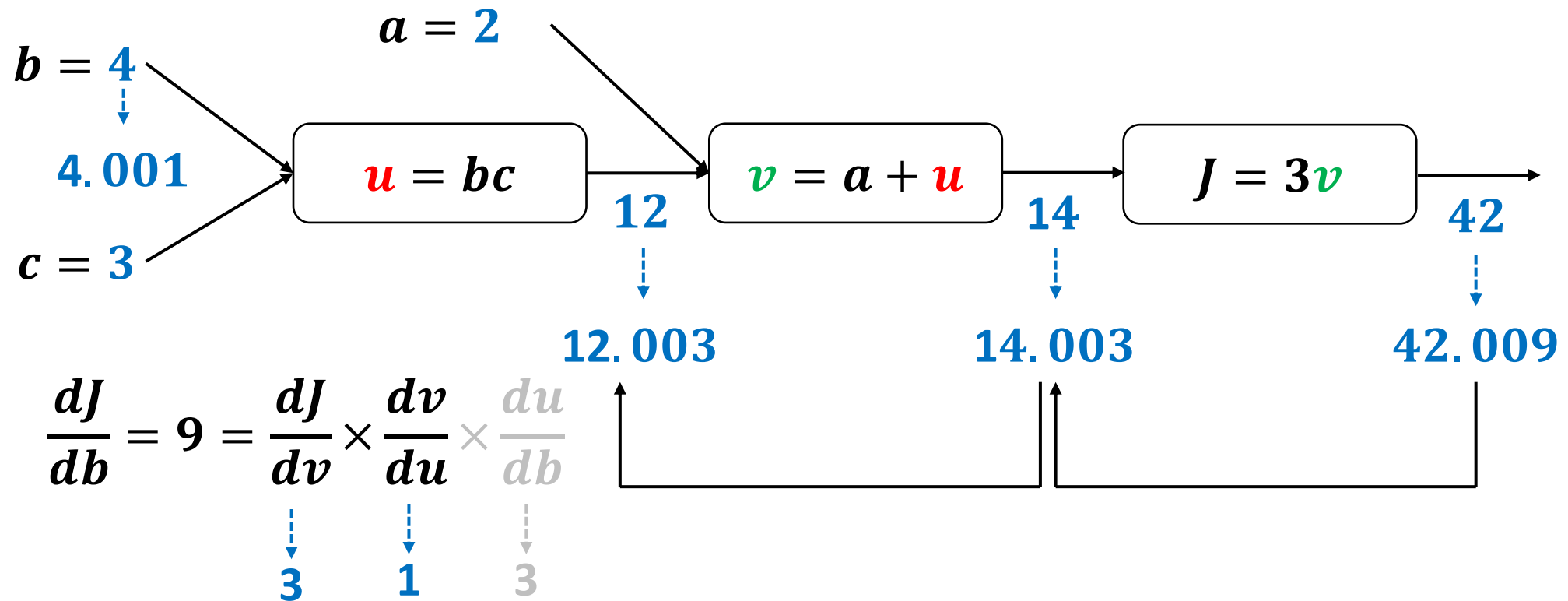
# Backward Propagation

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



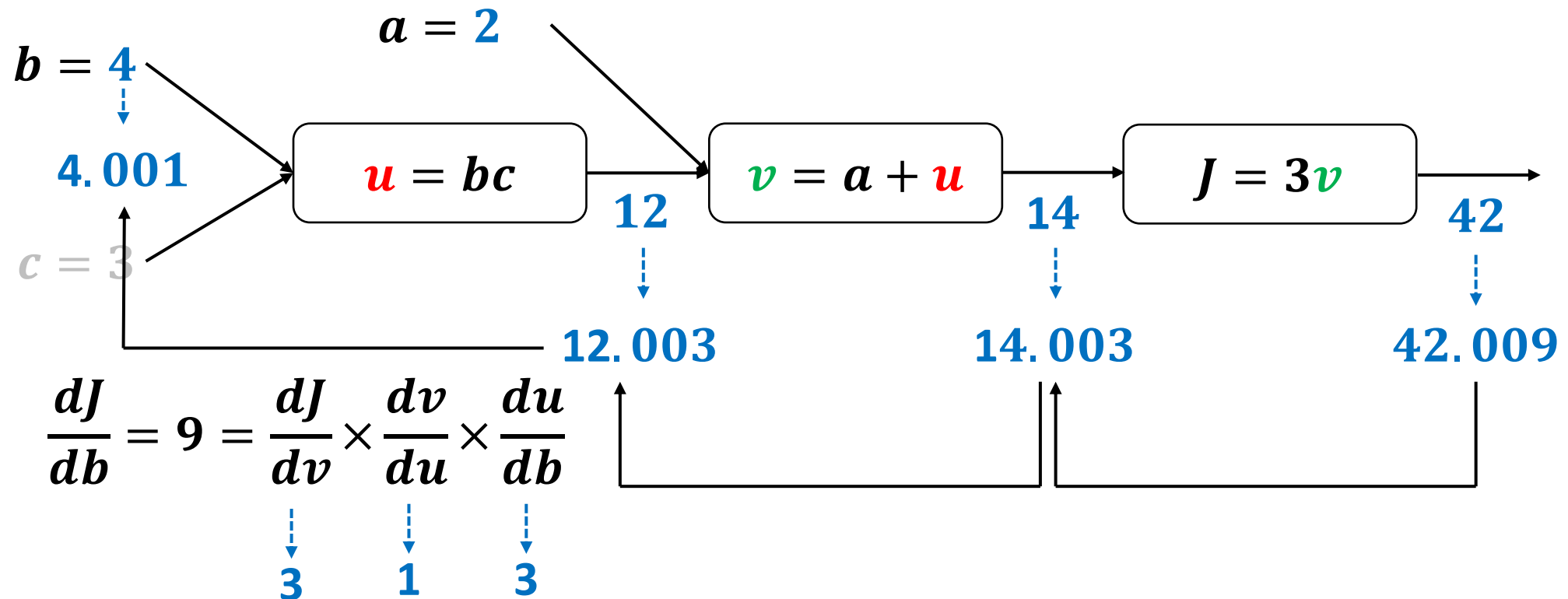
# Backward Propagation

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



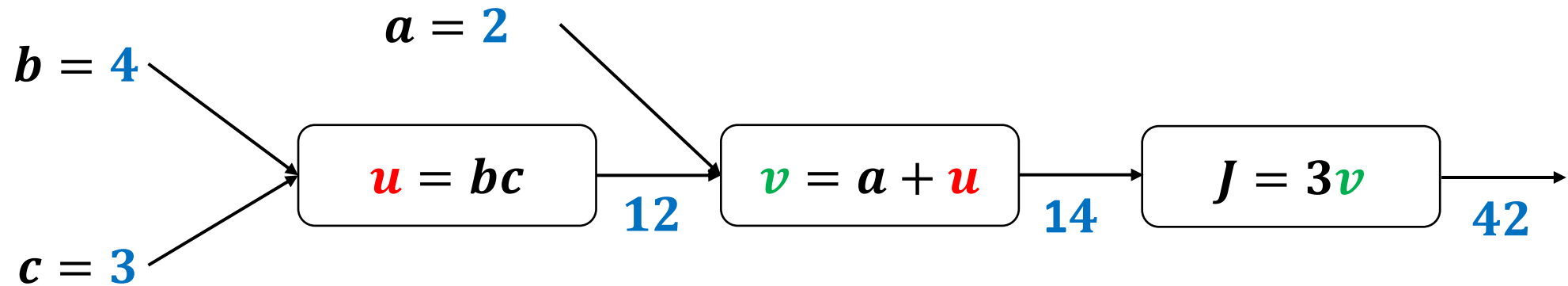
# Backward Propagation

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



# Backward Propagation

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

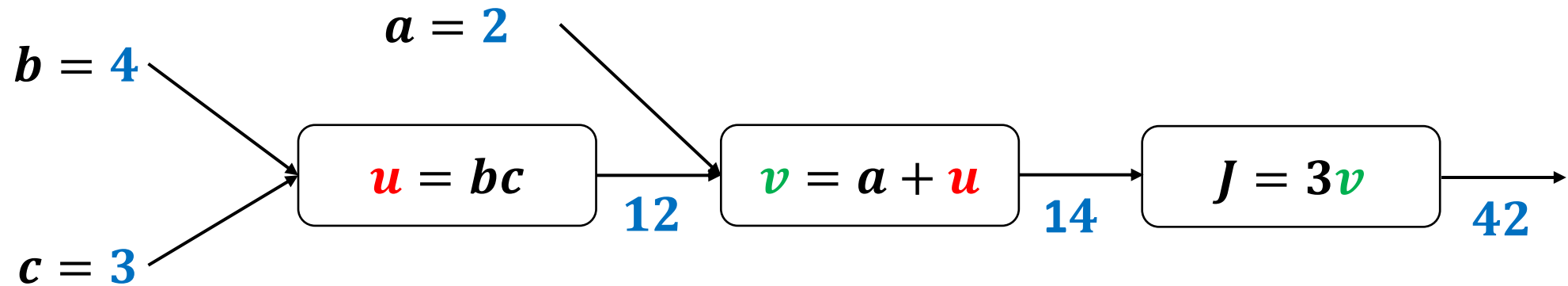


$$\frac{dJ}{dc} = \text{Derivative of } J \text{ with respect to } c$$



# Backward Propagation

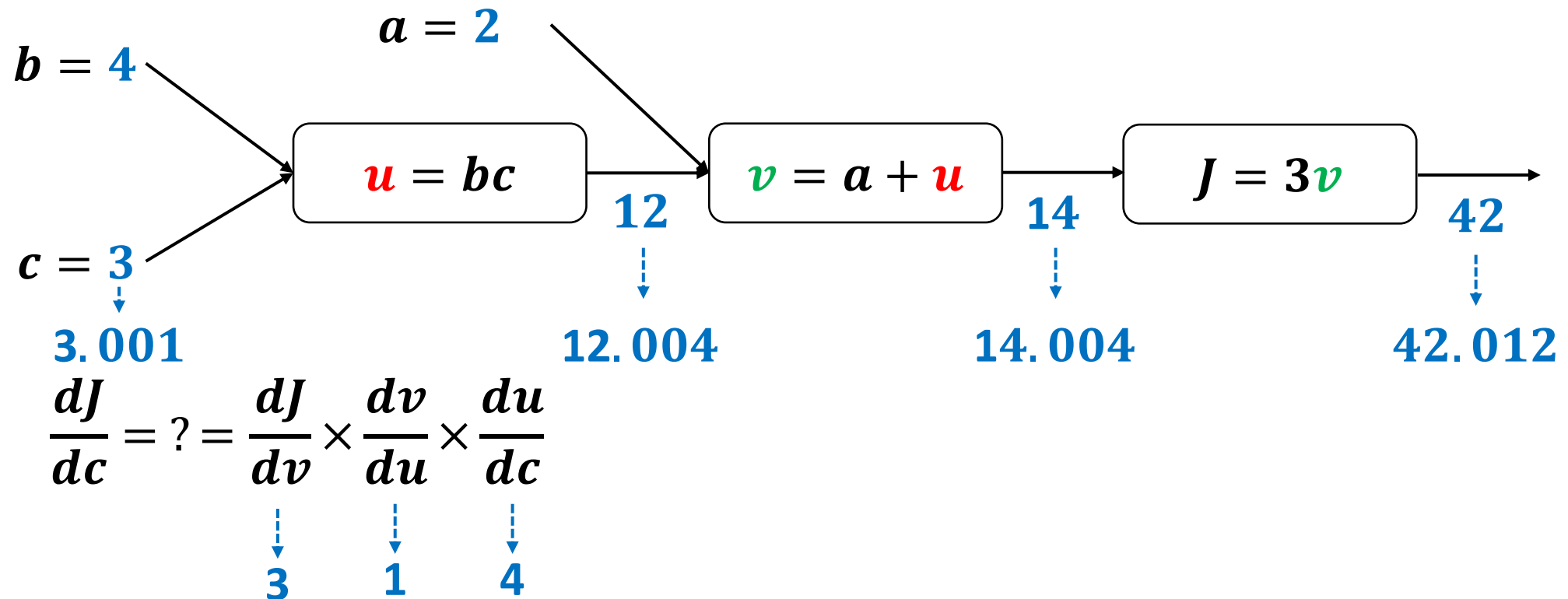
- Let us now compute the derivatives of the variables through the computation graph as follows:



$\frac{dJ}{dc}$  = If we change  $c$  a little bit, how would  $J$  change?

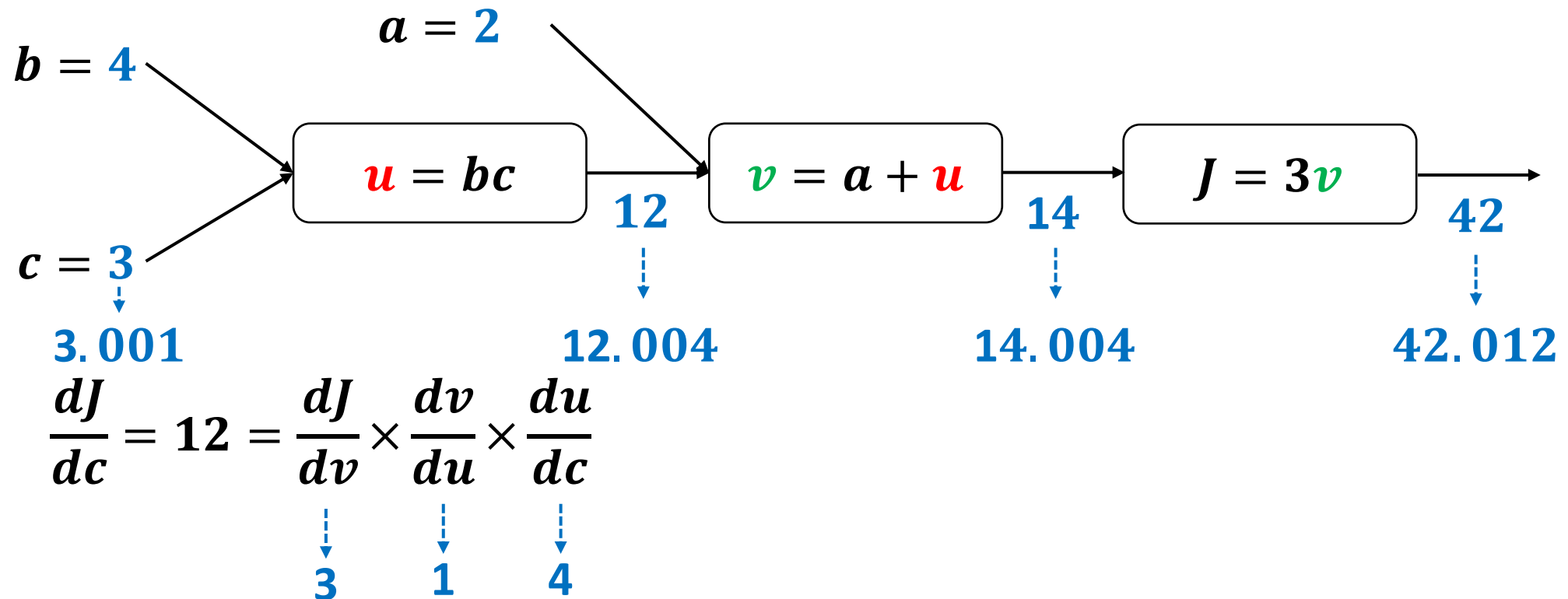
# Backward Propagation

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



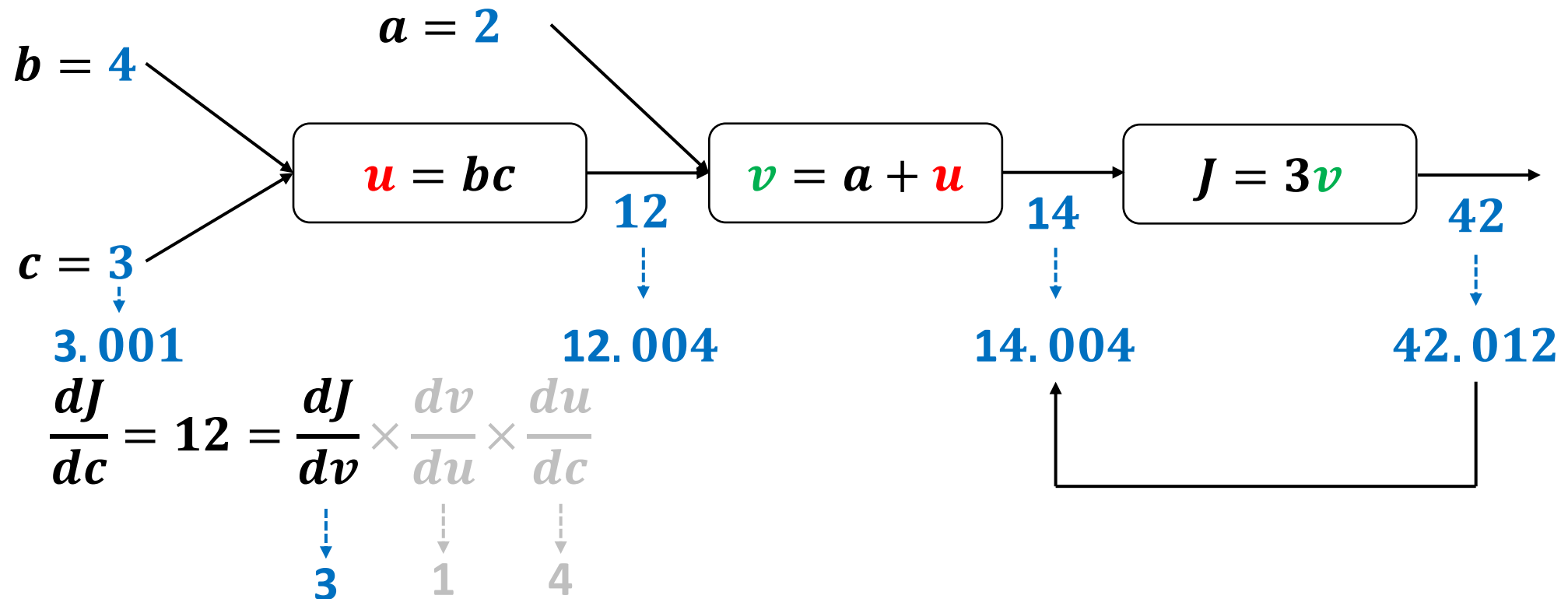
# Backward Propagation

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



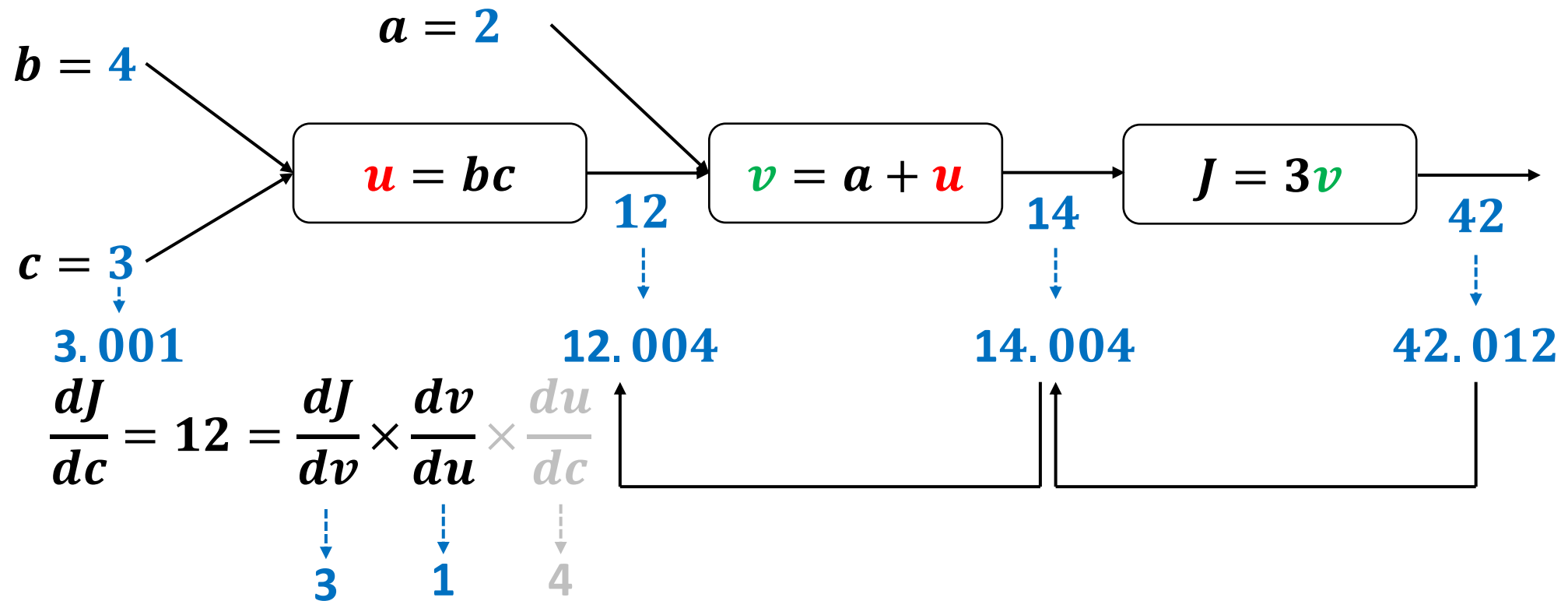
# Backward Propagation

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



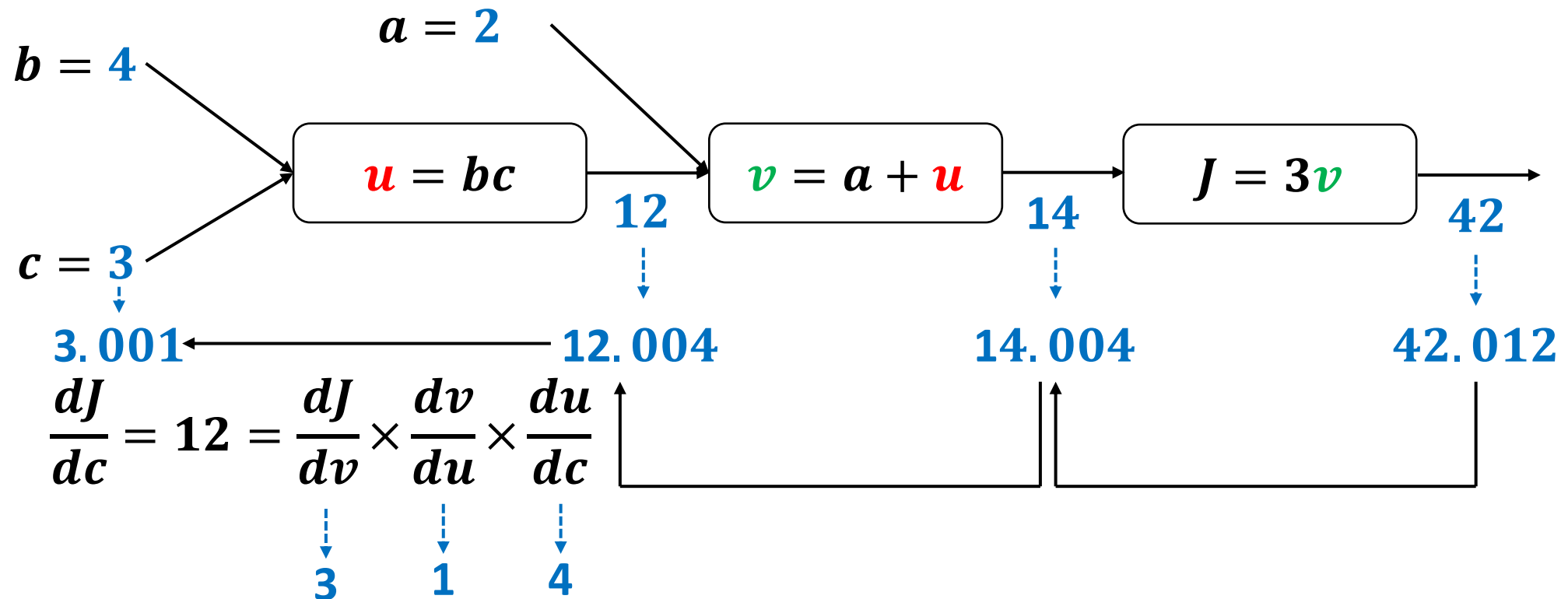
# Backward Propagation

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



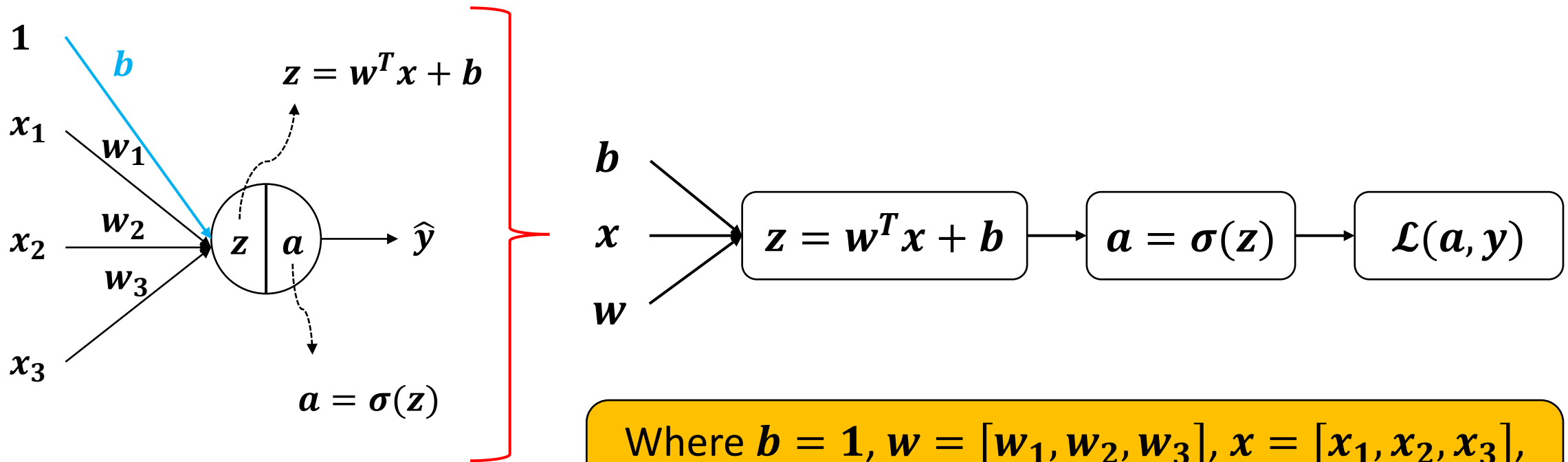
# Backward Propagation

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



# The Computation Graph of Logistic Regression

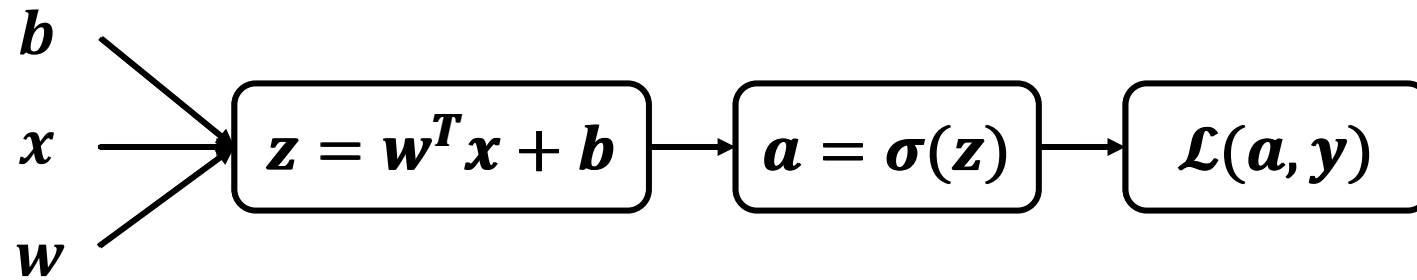
- Let us translate logistic regression (which is a **neural network** with only 1 neuron) into a computation graph



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

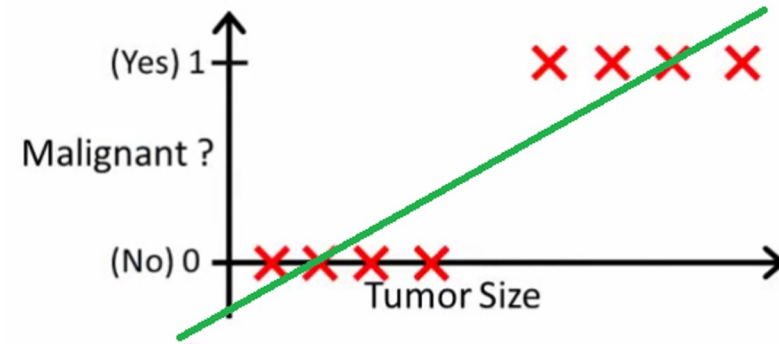




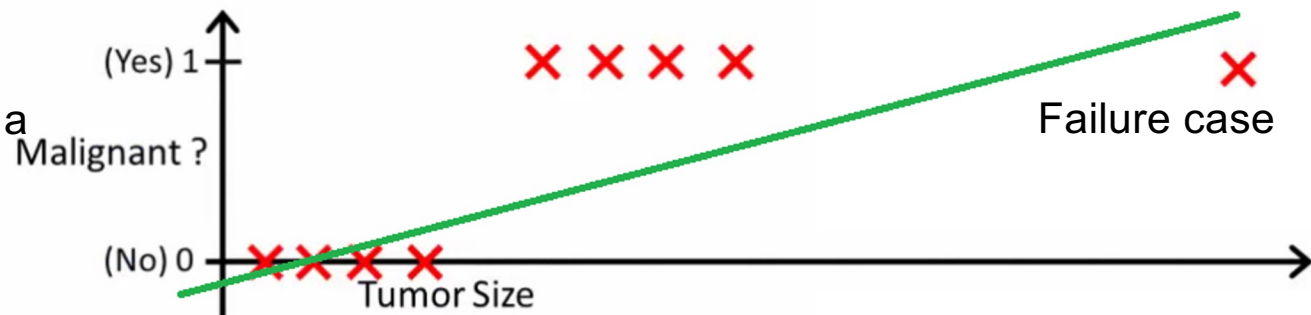
# Recap: Using Linear Regression for Binary Classification

## Example

- Using linear regression, fit a polynomial (line in this case) through the training data of type {tumor size, tumor type}.
- Malignant tumors means 1 and non-malignant means 0.
- Intuition: All tumors larger certain threshold are malignant
- Green line is the model  $h(x)$ .
- Prediction Rule: If  $h(x) > 0.5$ , for any given tumor size  $x$ , predict malignant tumor and benign otherwise.



Linear regression gives a raw number but logistic regression tells probability that  $x$  belongs to the "positive" class. Hence, it is a regression algorithm but, by setting a rule on the probability, we can perform classification.

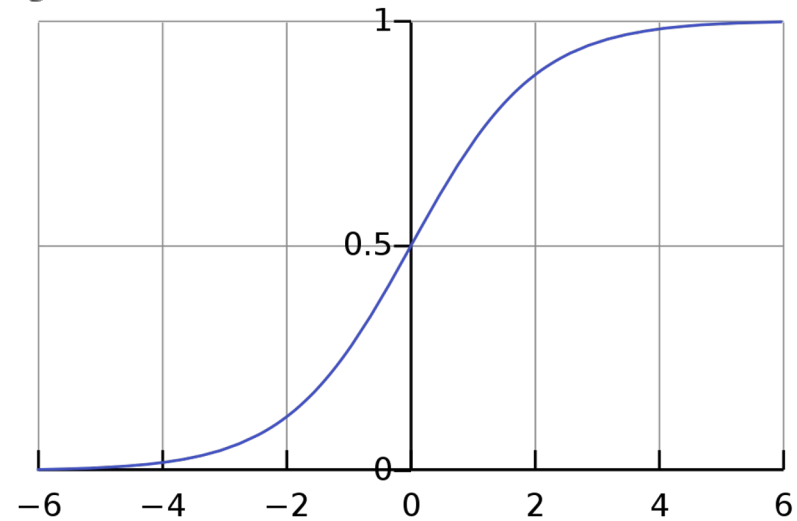


# Using Linear Regression for Binary Classification

*Turning scores of input  $\mathbf{x}$ ,  $\mathbf{w}^T \mathbf{x}$  into probabilities using sigmoid function*

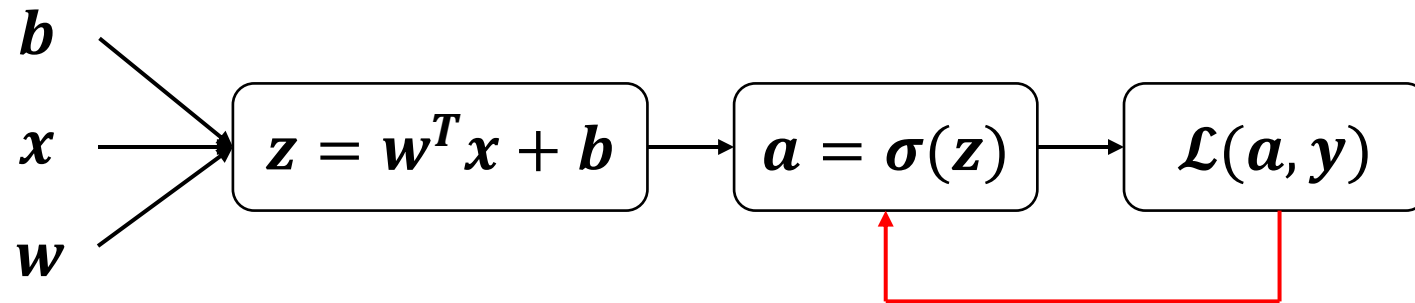
$$p(y = 1 | \mathbf{x}, \mathbf{w}) = a = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{e^{\mathbf{w}^T \mathbf{x}}}{1 + e^{\mathbf{w}^T \mathbf{x}}} = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

$$p(y = 0 | \mathbf{x}, \mathbf{w}) = 1 - a = 1 - \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}}}$$



# Backward Propagation

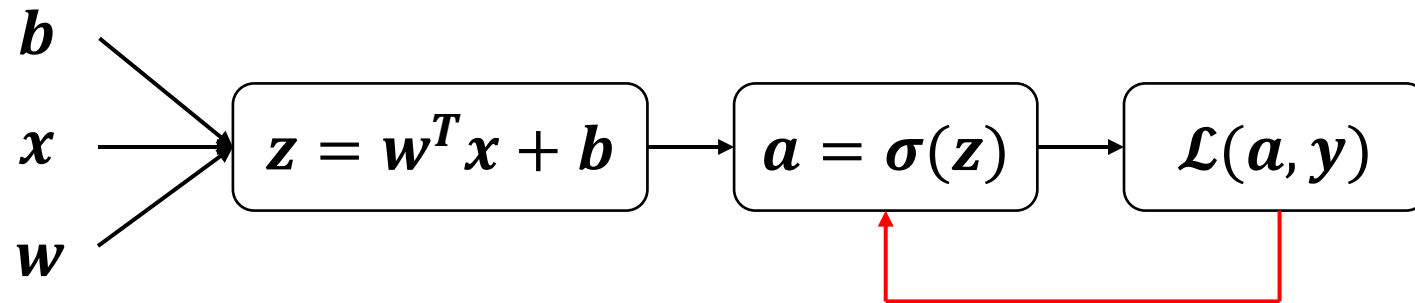
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

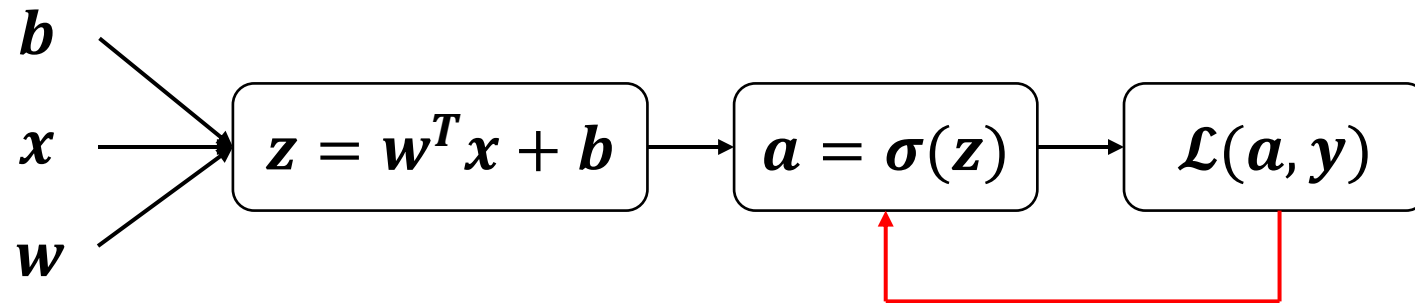
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

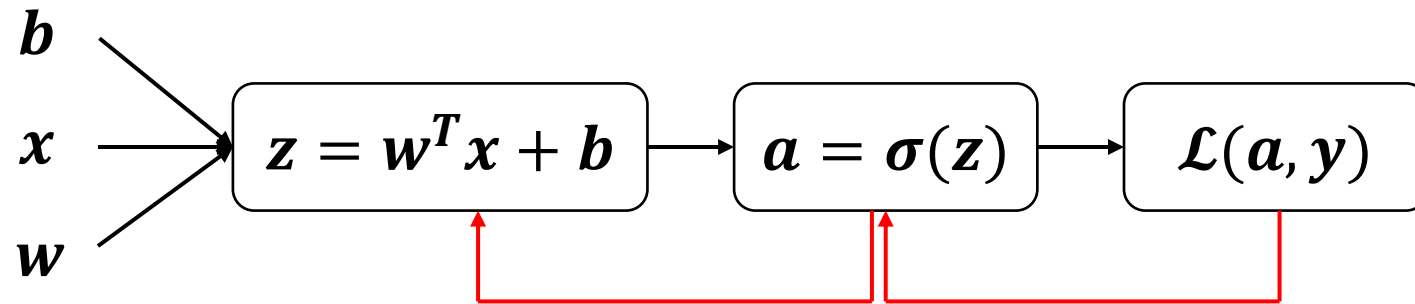
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

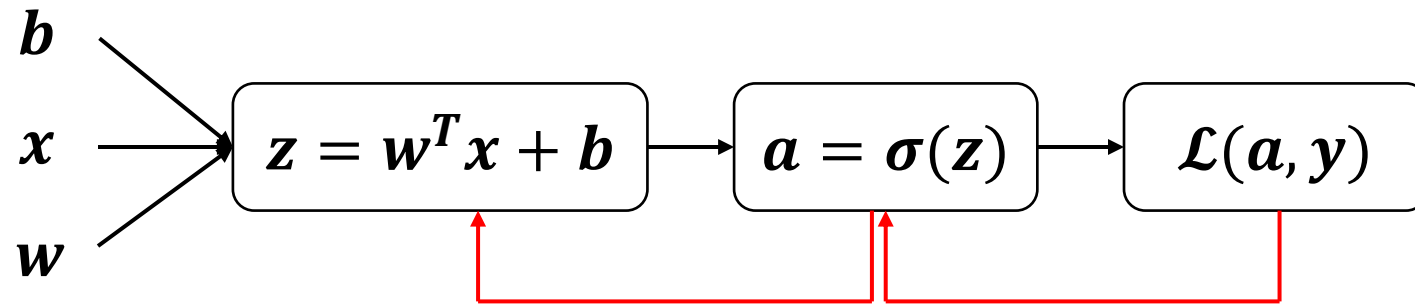
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

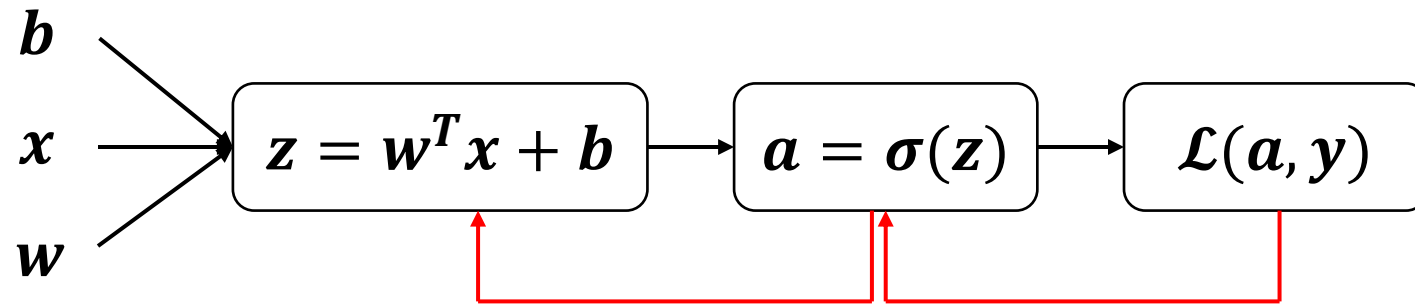
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

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

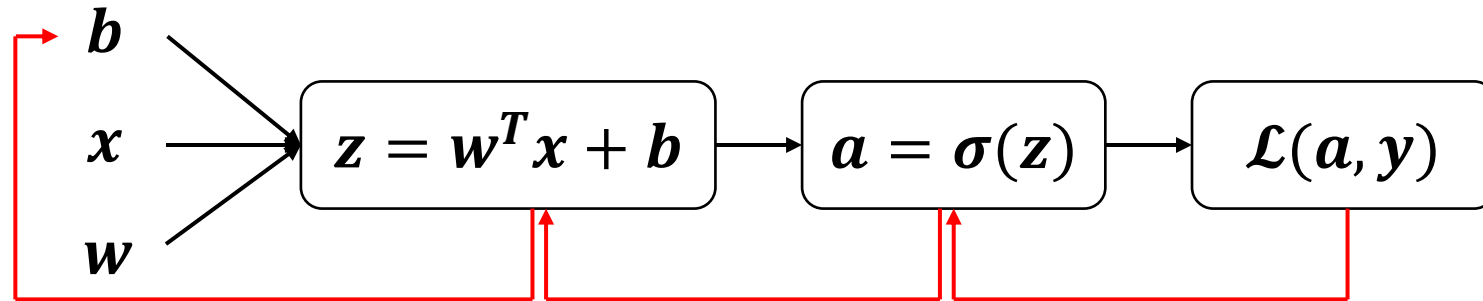


$$\frac{\partial \mathcal{L}}{\partial z} = a - y$$



# Backward Propagation

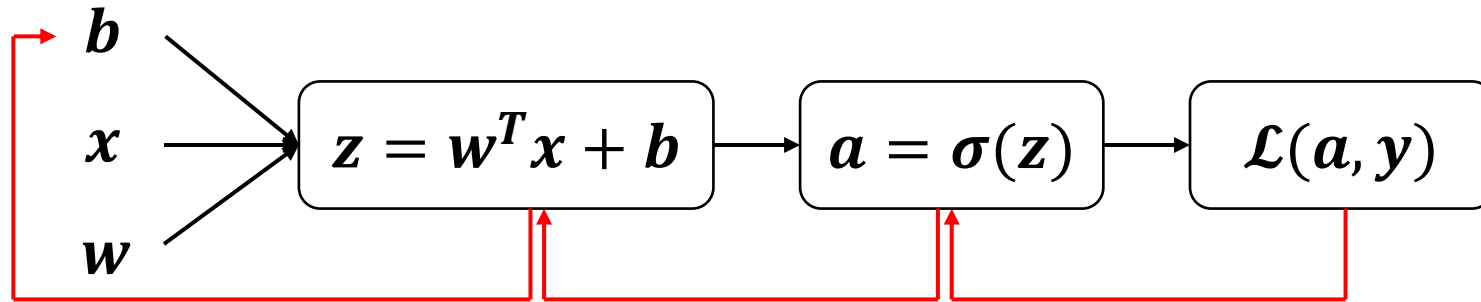
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

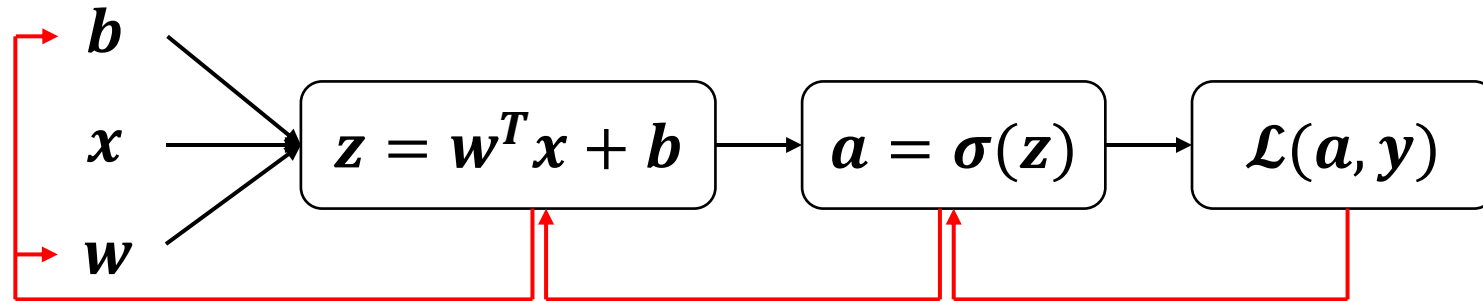
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

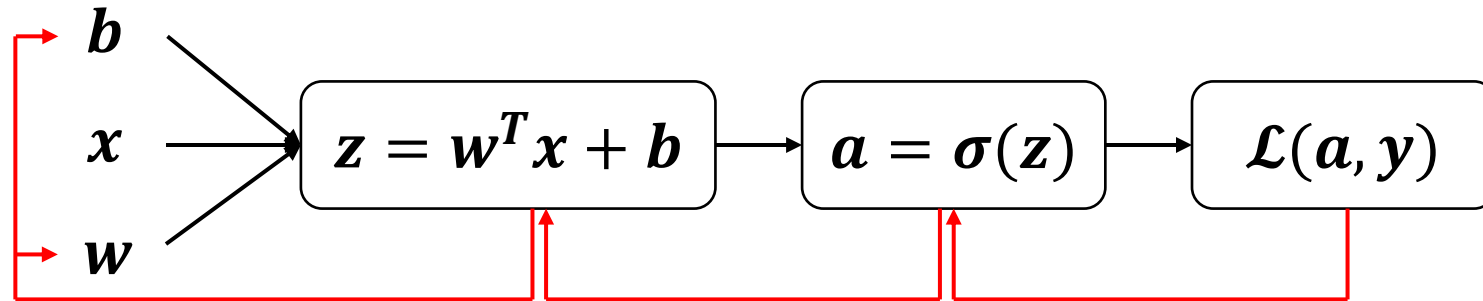
- The derivatives can be computed by moving from right to left



$\frac{\partial \mathcal{L}}{\partial w}$  = Partial derivative of  $\mathcal{L}$  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:

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


We will denote this as  $\textcolor{blue}{dz}$  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  $\mathcal{L}(\mathbf{w}, \mathbf{b})$ , where  $\mathbf{w} = [\mathbf{w}_1, \dots, \mathbf{w}_m]$  and  $\mathbf{b} = \mathbf{w}_0$
- Start off with some guesses for  $\mathbf{w}_1, \dots, \mathbf{w}_m$ 
  - It does not really matter what values you start off with for  $\mathbf{w}_1, \dots, \mathbf{w}_m$ , but a common choice is to set them all initially to zero

- Repeat until convergence{

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

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

}

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

$$\begin{cases} z^{(i)} = w^T x^{(i)} + b \\ a^{(i)} = \sigma(z^{(i)}) \end{cases}$$

Backward  
propagation

$$\begin{cases} dz^{(i)} = a^{(i)} - y^{(i)} \\ dw = dw + dz^{(i)} x^{(i)} \\ db = db + dz^{(i)} \end{cases}$$

Outside  
the loop

$$\begin{cases} dw = dw/n \\ db = db/n \\ w = w - \alpha dw \\ b = b - \alpha db \end{cases}$$

}

$$Z = w^T X + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

$$dw = \frac{1}{n} X dZ^T$$

$$db = \frac{1}{n} \sum_{i=1}^n dz^{(i)}$$

$$w = w - \alpha dw$$

$$b = b - \alpha db$$

Vectorized version



# Gradient Descent For Logistic Regression

- **Outline:**

- Repeat until convergence{

$$\mathbf{Z} = \mathbf{w}^T \mathbf{X} + \mathbf{b}$$

$$\mathbf{A} = \sigma(\mathbf{Z})$$

$$d\mathbf{Z} = \mathbf{A} - \mathbf{Y}$$

$$d\mathbf{w} = \frac{1}{n} \mathbf{X} d\mathbf{Z}^T$$

$$d\mathbf{b} = \frac{1}{n} \sum_{i=1}^n dz^{(i)}$$

$$\mathbf{w} = \mathbf{w} - \alpha d\mathbf{w}$$

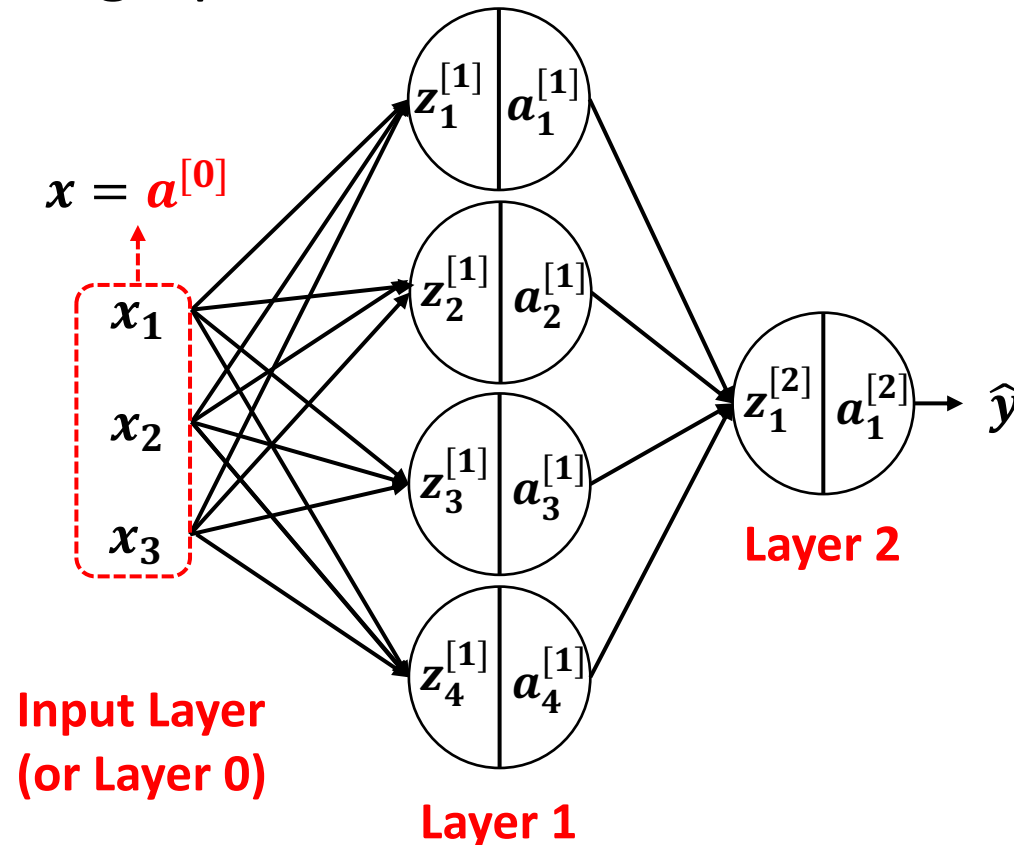
$$\mathbf{b} = \mathbf{b} - \alpha d\mathbf{b}$$

}

# The Computation Graph of A Neural Network

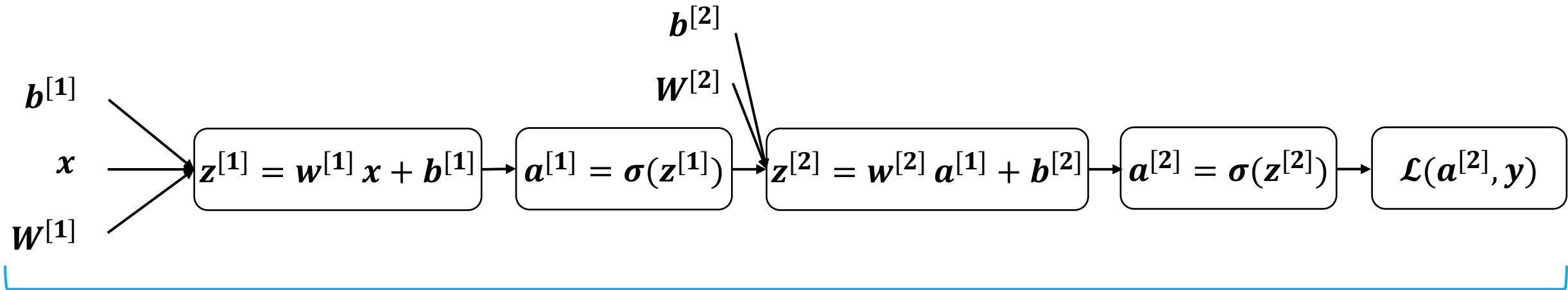
- Akin to logistic regression, we can represent any neural network in terms of a computation graph

A neural network  
with 2 layers



# The Computation Graph of A Neural Network

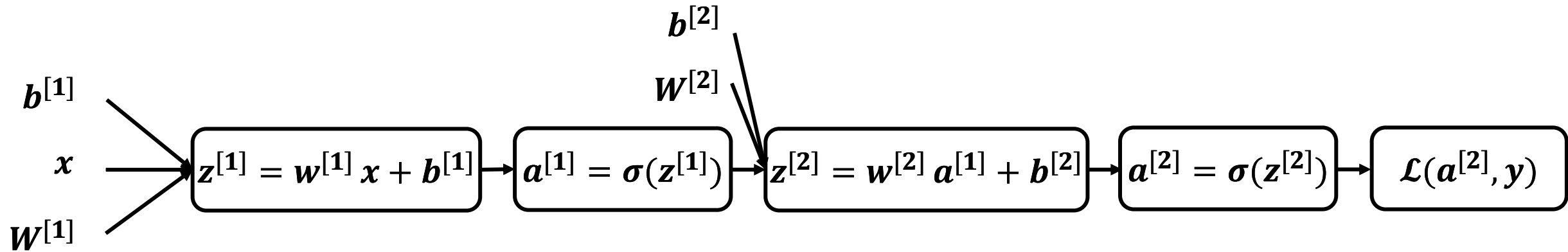
- Akin to logistic regression, we can represent any neural network in terms of a computation graph



The corresponding computation graph

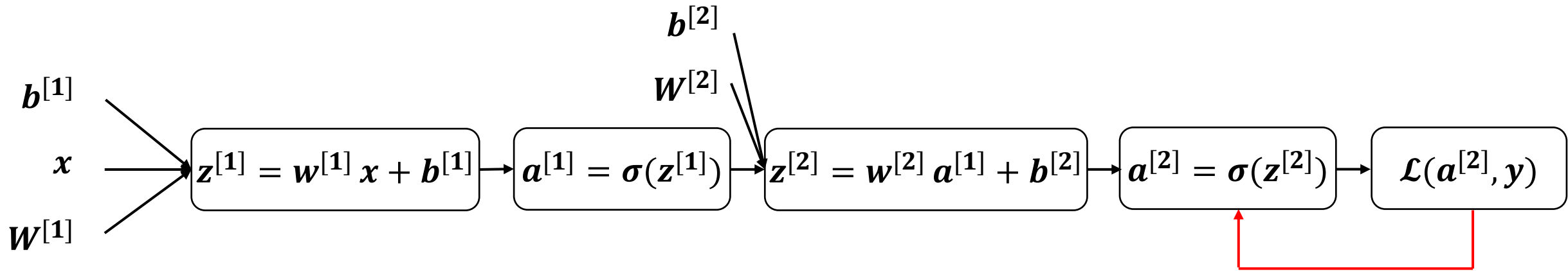
# Forward Propagation

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



# Backward Propagation

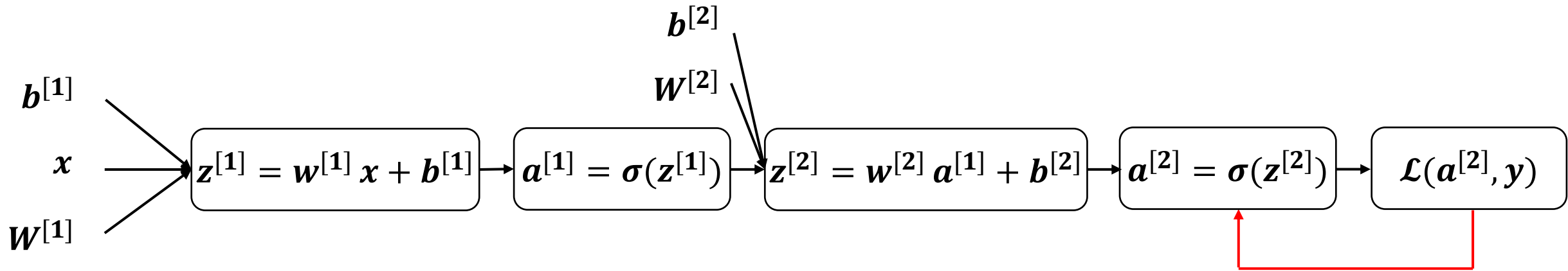
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

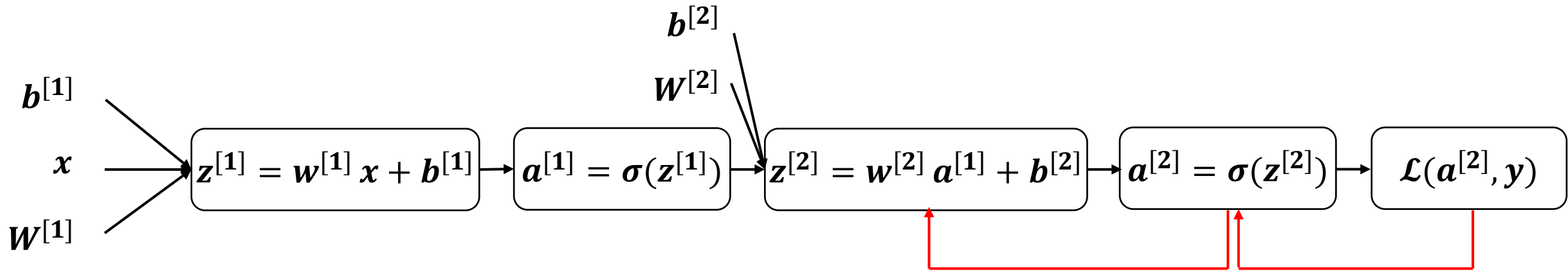
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

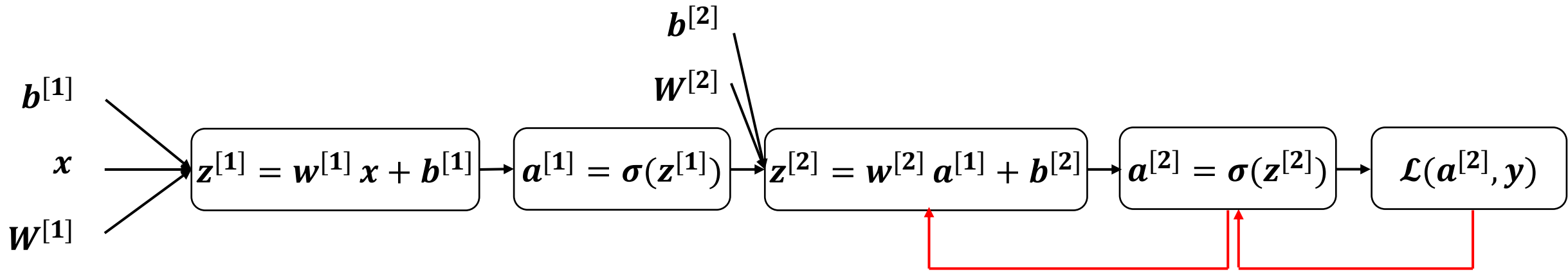
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

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

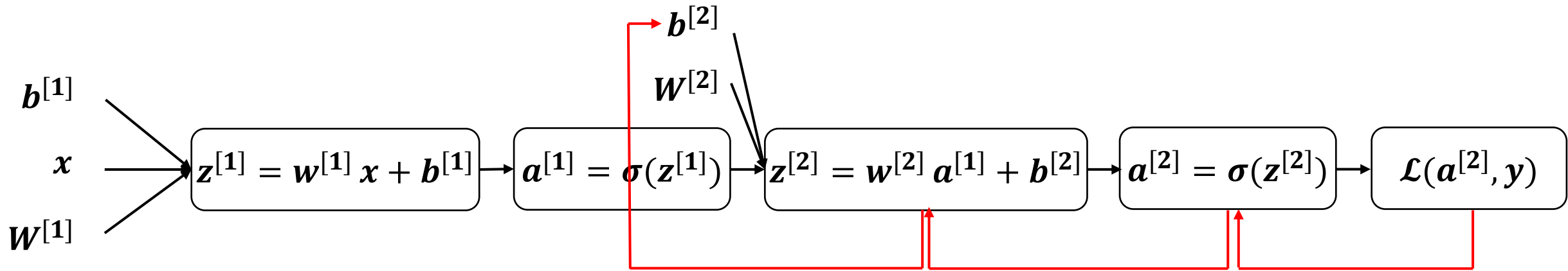


$$\frac{\partial \mathcal{L}}{\partial z^{[2]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} = a^{[2]} - y$$



# Backward Propagation

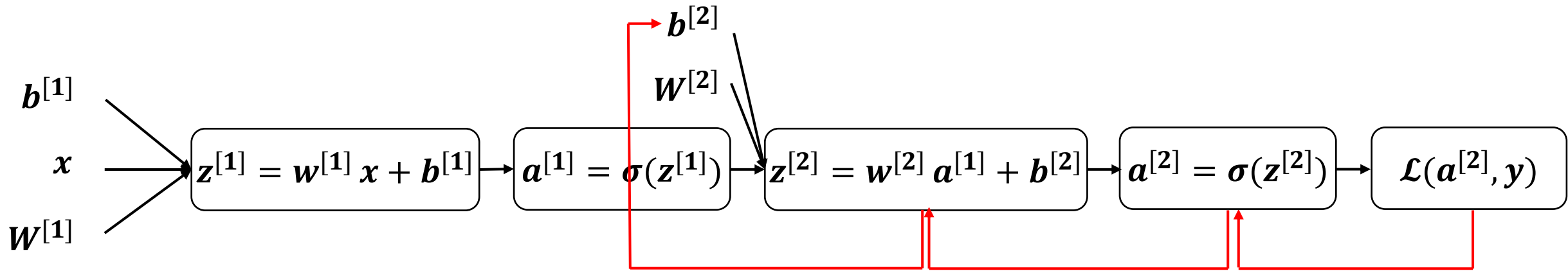
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

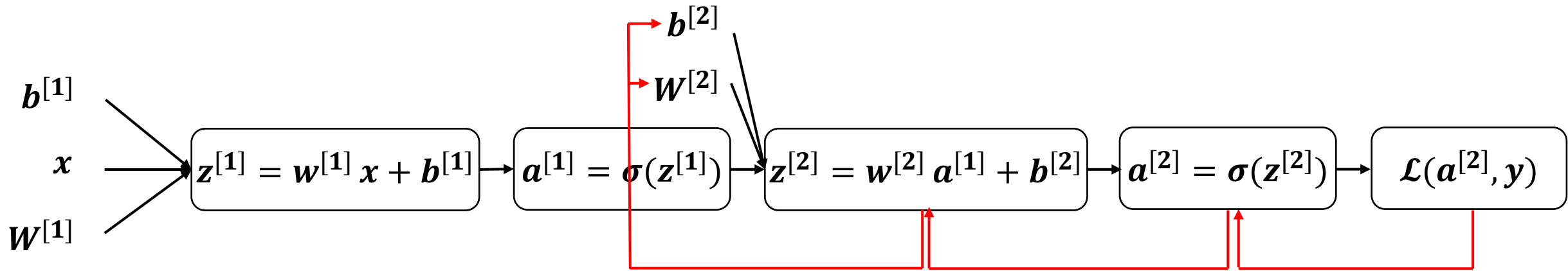
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

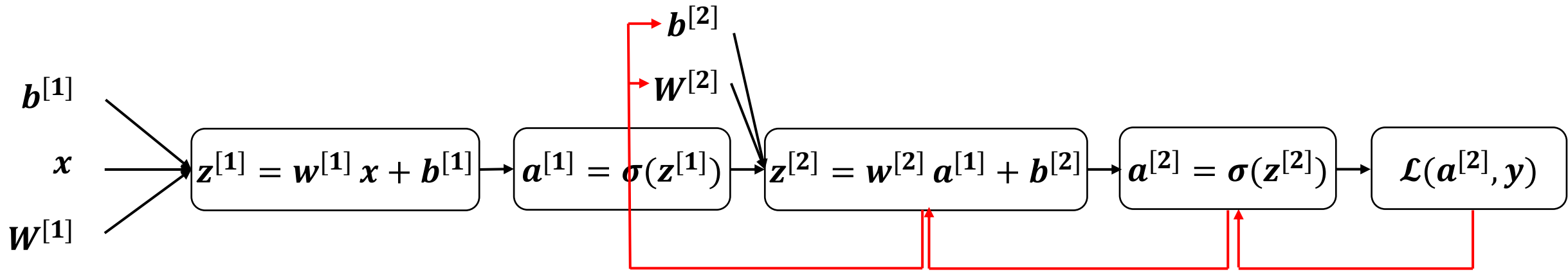
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

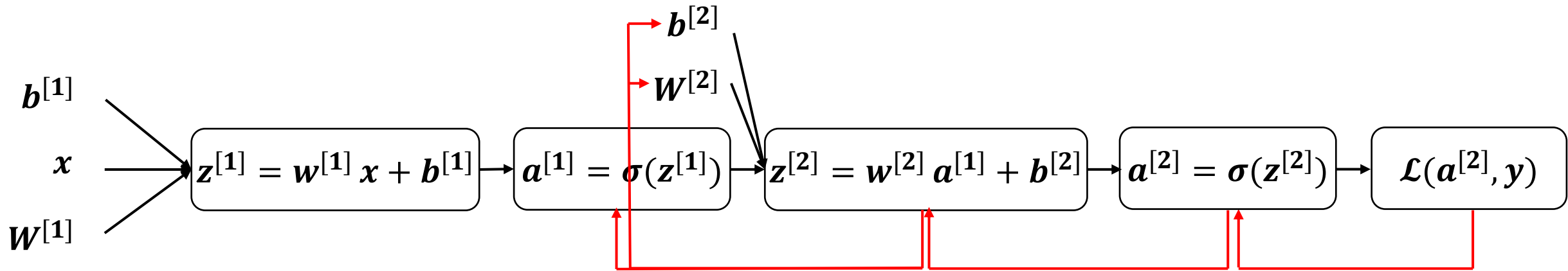
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

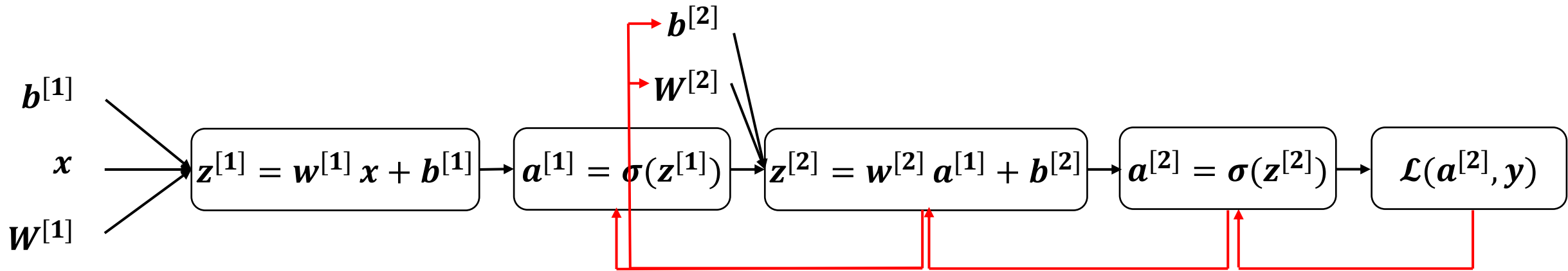
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

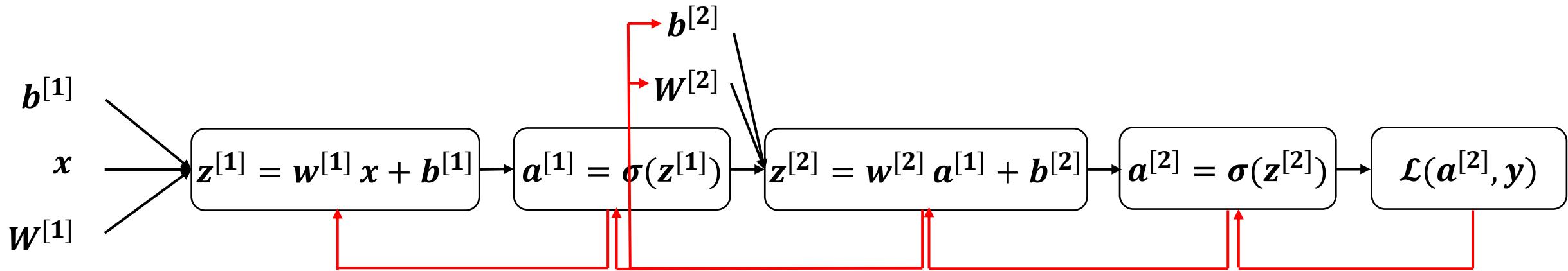
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

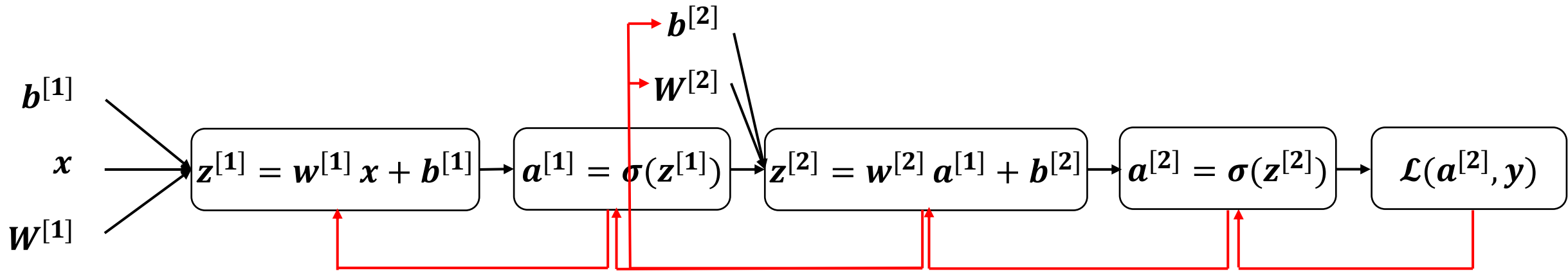
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

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



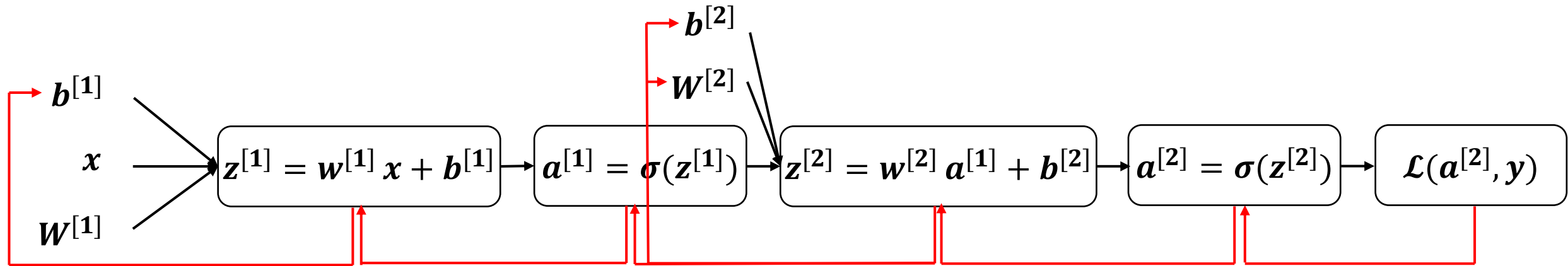
$$\frac{\partial \mathcal{L}}{\partial z^{[1]}} = \frac{\partial \mathcal{L}}{\partial a^{[2]}} \times \frac{\partial a^{[2]}}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial a^{[1]}} \times \frac{\partial a^{[1]}}{\partial z^{[1]}} = (a^{[2]} - y)w^{[2]T} * a^{[1]}(1 - a^{[1]})$$

Element-wise product



# Backward Propagation

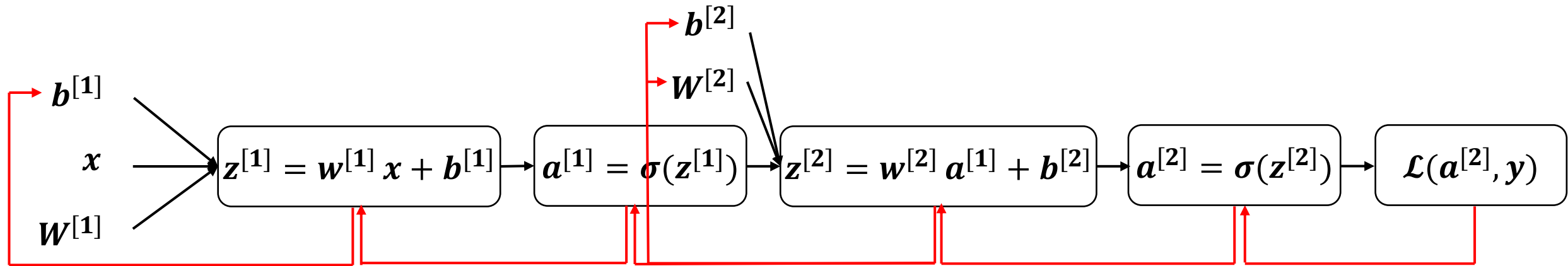
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

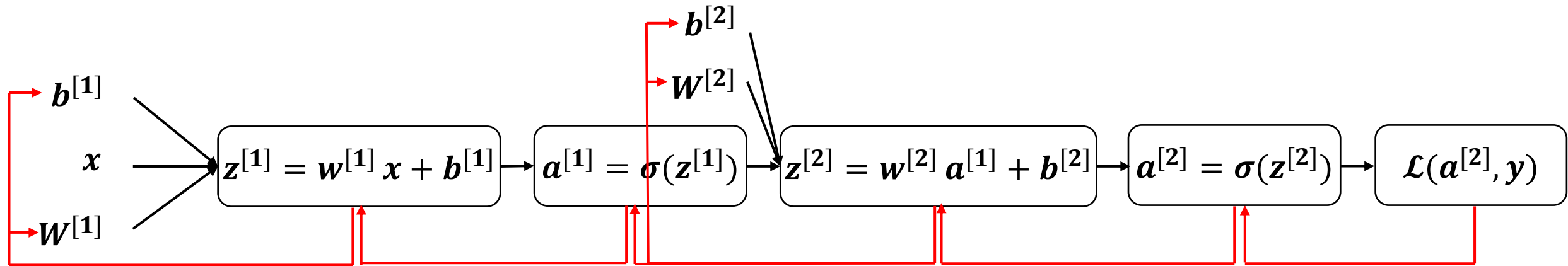
- The derivatives can be computed by moving from right to left



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

# Backward Propagation

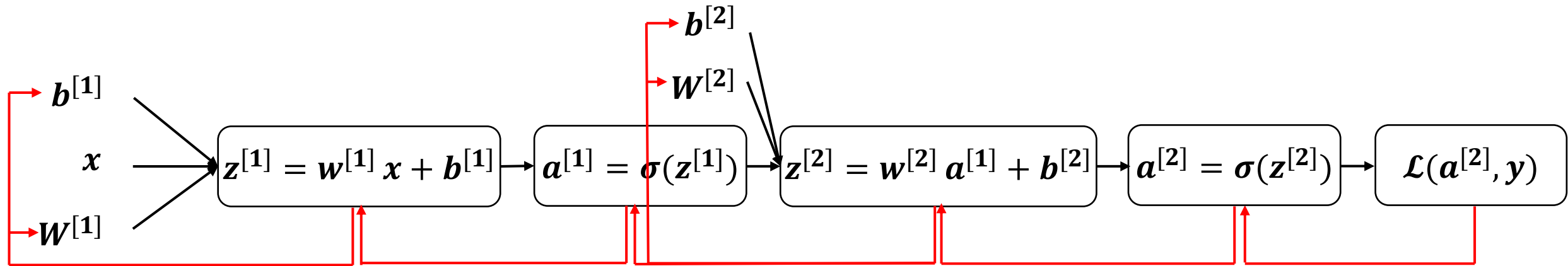
- The derivatives can be computed by moving from right to left



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

# 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:

$$dz^{[2]} = \frac{\partial \mathcal{L}}{\partial z^{[2]}} = a^{[2]} - y$$

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

$$db^{[2]} = \frac{\partial \mathcal{L}}{\partial b^{[2]}} = a^{[2]} - y$$

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

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

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