

# Sixteenth Session

**Alireza Moradi**

---



| [Linkedin](#)



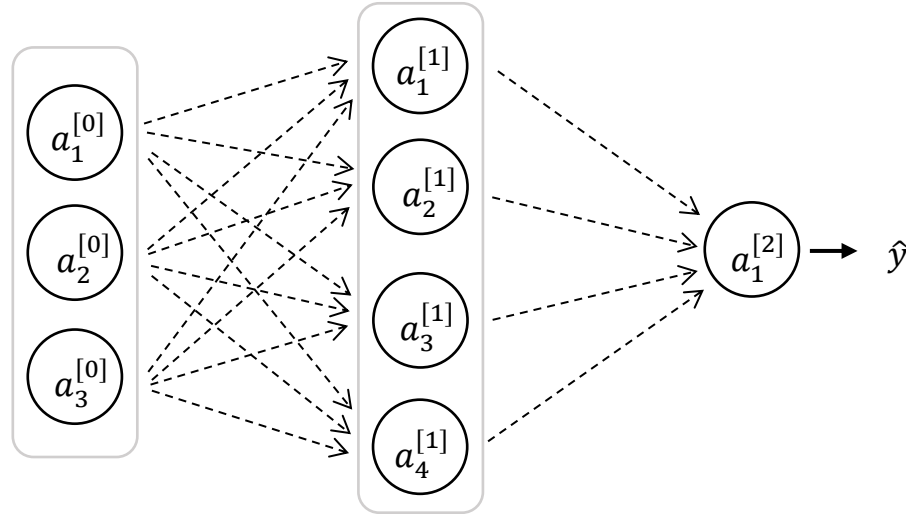
| [k](#)

# **Backpropagation algorithm**

---

# ANN learning process

- Suppose we want to use the following Artificial Neural Network (ANN) to predict a disease:



- Before using this architecture for prediction, we must **learn the parameters**.
- ✓ In this example, how many parameters do we have?

# ANN learning process

$$\# \text{ Parameters} = \# \text{ weights} + \# \text{ biases}$$

- **Backpropagation** algorithm is the standard method for **learning weights and biases**.
- **Backpropagation** enables the network (ANN) to learn by determining the exact **changes to make to weights and biases** to produce an accurate result.

**Training process => Learning parameters => Backpropagation**

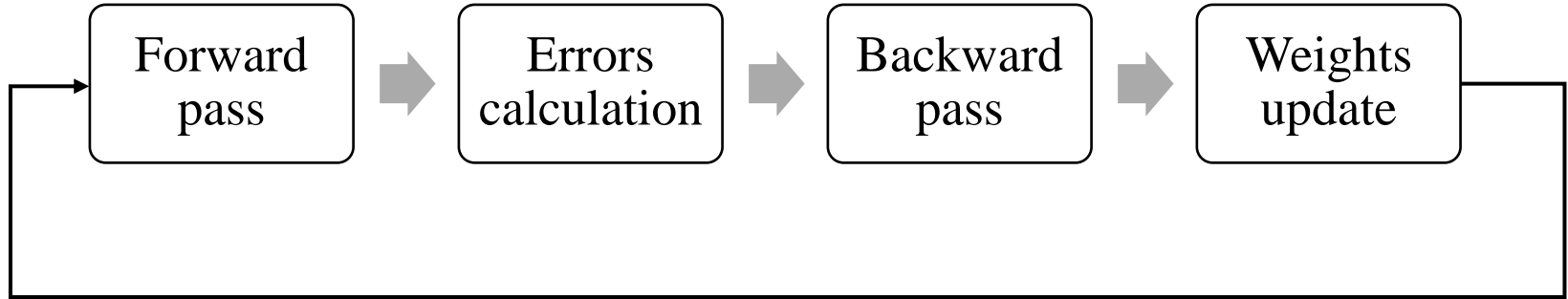
- We must **initialize the network parameters** to something reasonable and use **Backpropagation** algorithm to adjust them.

# Backpropagation Algorithm

- You can think of **Backpropagation Algorithm** as a **feedback system** where, after each round of training, the network **reviews its performance** on tasks.
  - For this purpose, it calculates the **difference** between its output and the correct answer, known as the **error**.
  - Then, it **adjusts its internal parameters** to reduce this error next time.
  - Backpropagation algorithm was introduced in the **1970s**.
- ✓ How does **Gradient Descent** help us learn the best parameters in linear regression?

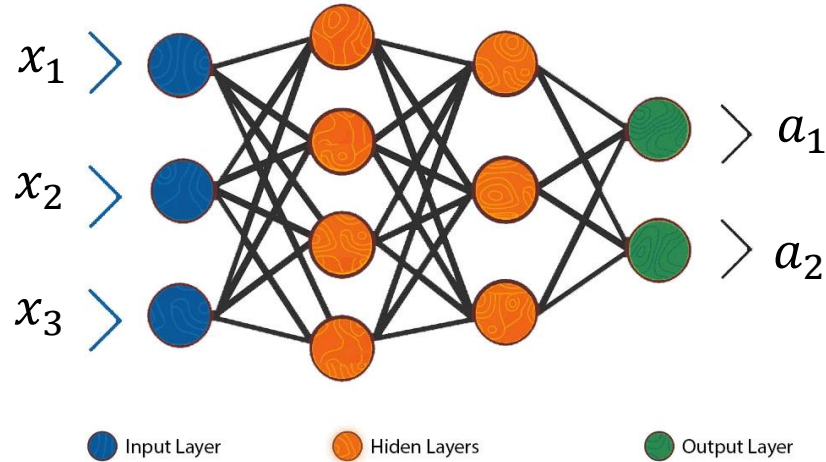
# Backpropagation Algorithm

- There are overall **four main steps** in the backpropagation algorithm:
  1. Forward pass
  2. Errors calculation
  3. Backward pass
  4. Parameters update



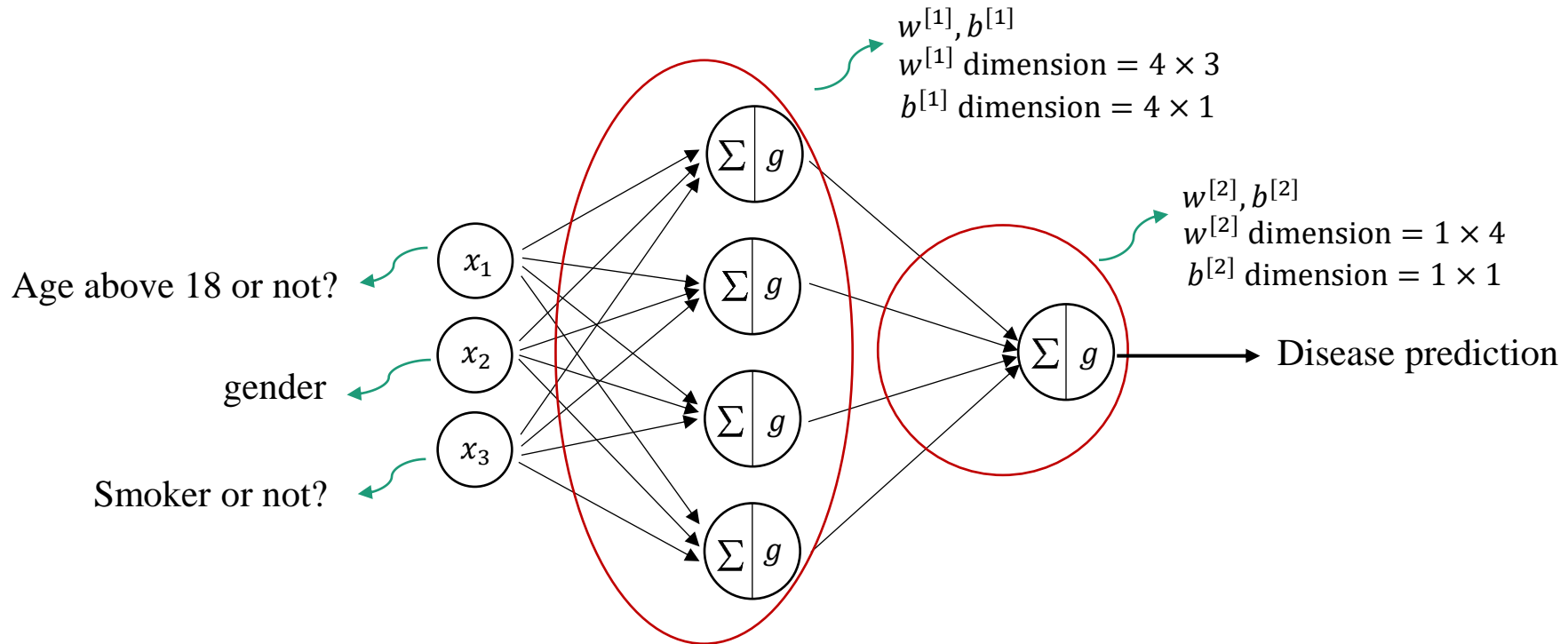
# Forward pass

- **Forward propagation**, also known as the **forward pass**, is the process of **computing the output of a neural network** given a set of input data.
- It involves **passing the input data through the network's layers**, applying weights and activation functions, and producing predictions or outputs.



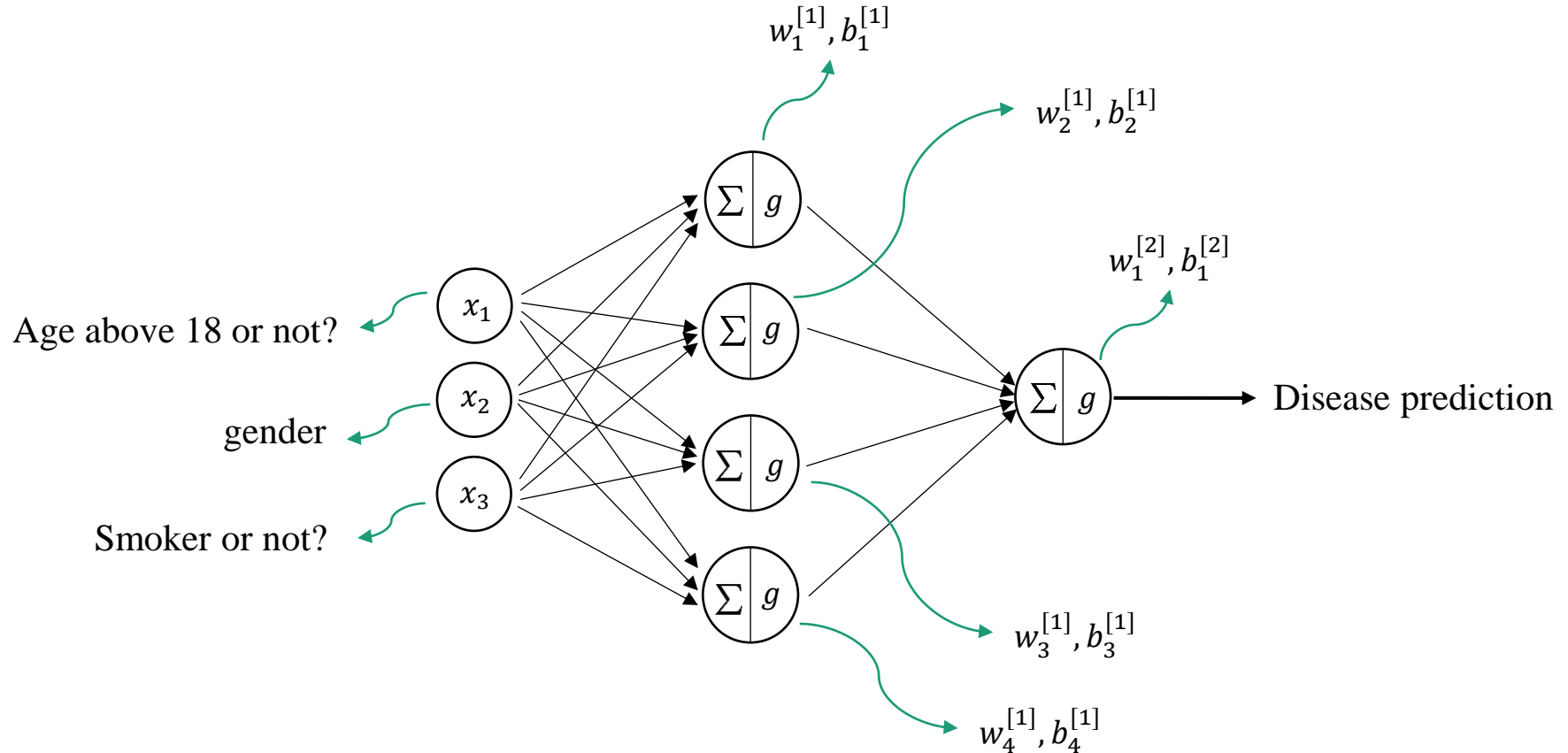
# Forward pass; Example

- We want to determine whether a person with certain features has diabetes or not?

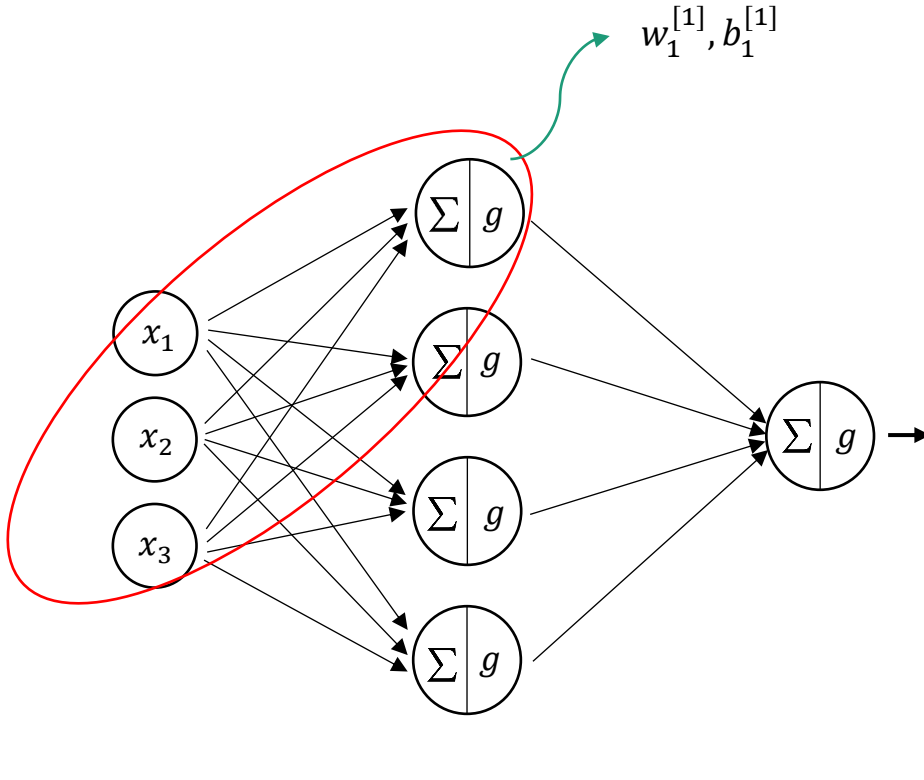




# Forward pass; Example



# Forward pass; Example



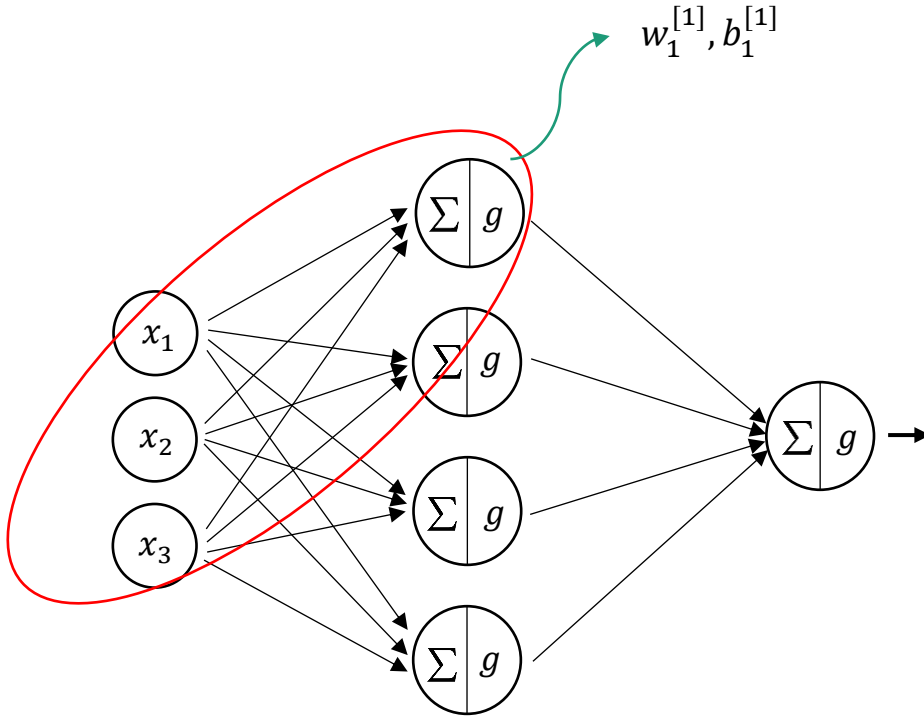
$$w_1^{[1]} = [0.1 \quad 0.2 \quad 0.5]$$

$$b_1^{[1]} = [0]$$

$$x = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$g = \text{sigmoid}$$

# Forward pass; Example



$$w_1^{[1]} = [0.1 \quad 0.2 \quad 0.5] \quad b_1^{[1]} = [0]$$

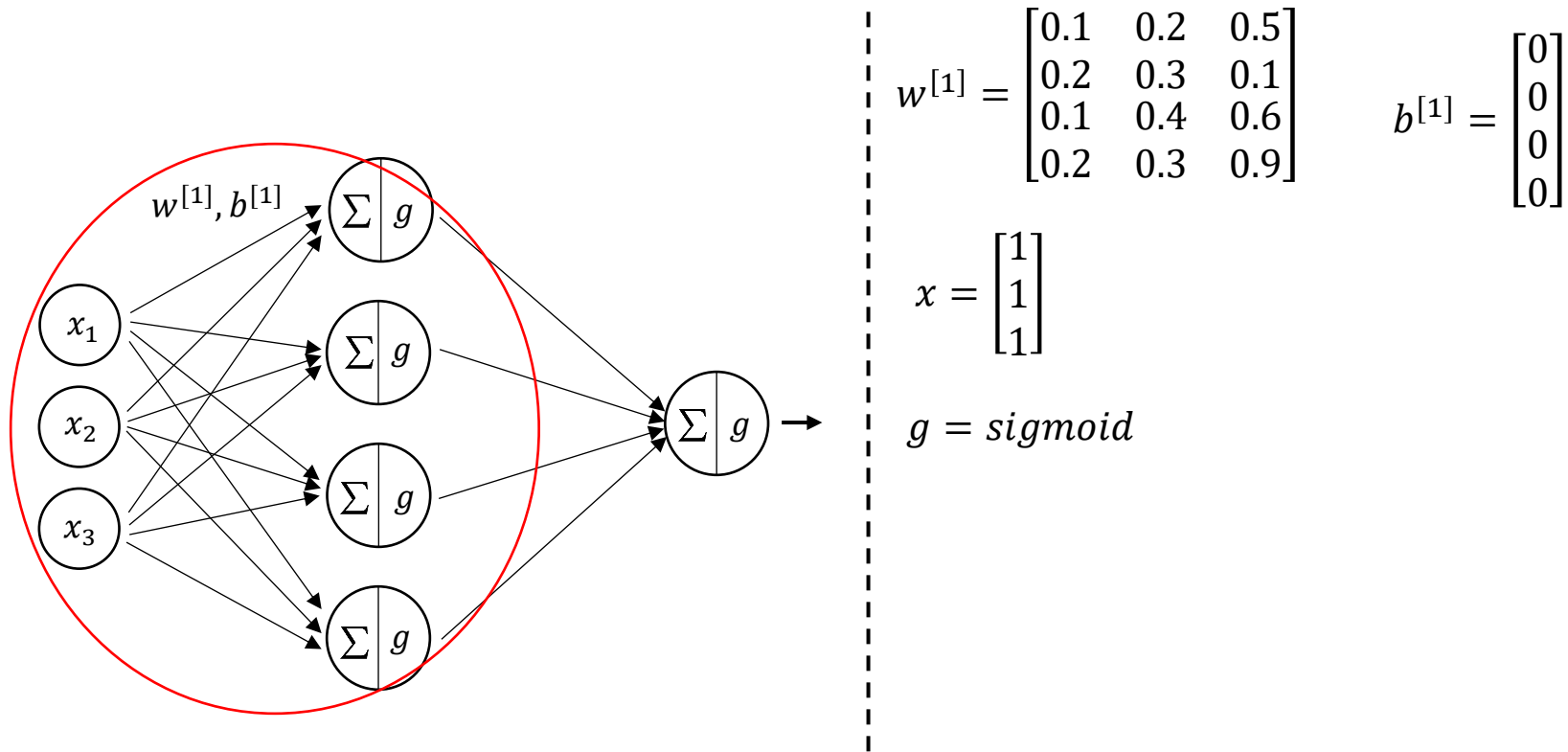
$$x = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$g = \text{sigmoid}$$

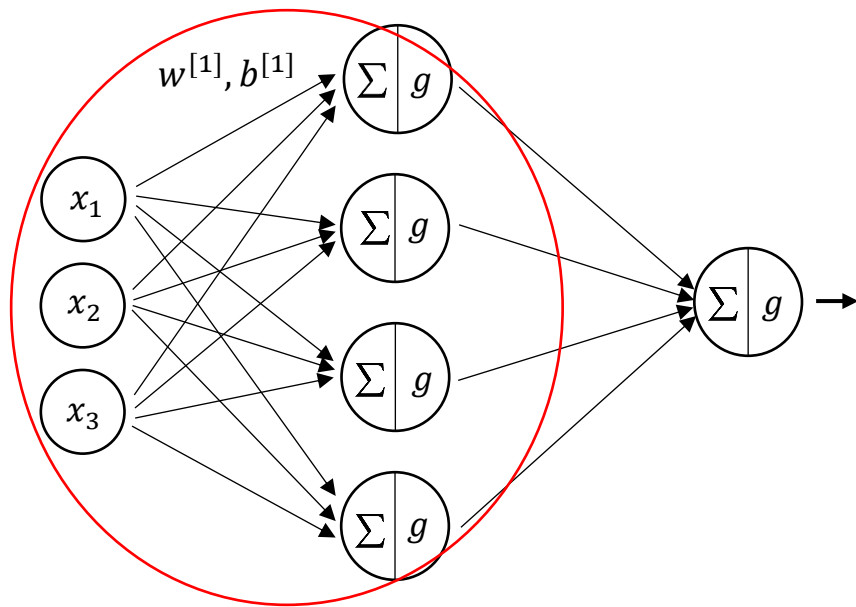
$$z_1^{[1]} = w_1^{[1]} \cdot x + b_1^{[1]} = \\ 0.1 \times 1 + 0.2 \times 1 + 0.5 \times 1 = 0.8$$

$$a_1^{[1]} = \text{sigmoid}(0.8) = \mathbf{0.68}$$

# Forward pass; Example



# Forward pass; Example



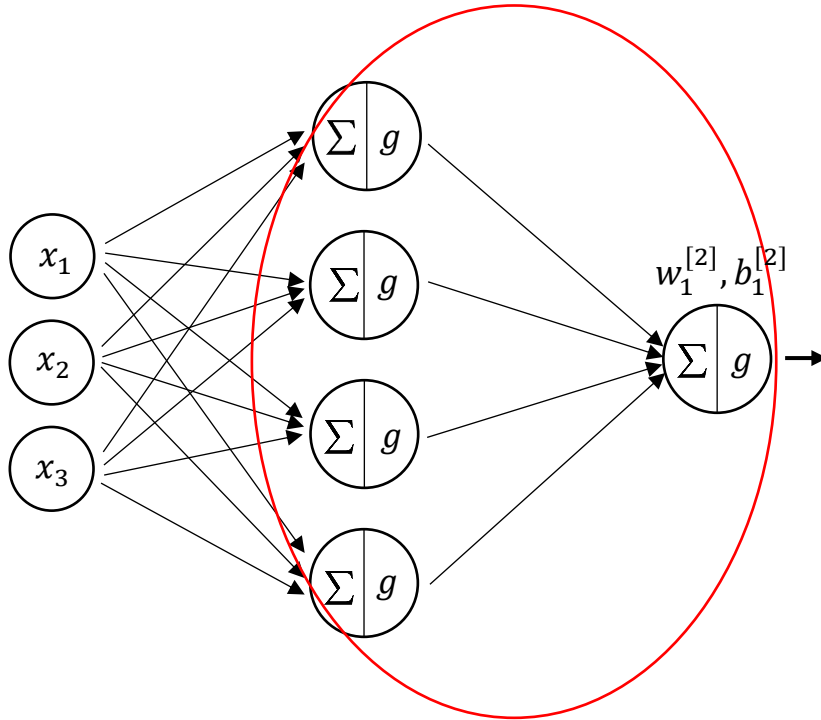
$$w^{[1]} = \begin{bmatrix} 0.1 & 0.2 & 0.5 \\ 0.2 & 0.3 & 0.1 \\ 0.1 & 0.4 & 0.6 \\ 0.2 & 0.3 & 0.9 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$x = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad g = \text{sigmoid}$$

$$z^{[1]} = w^{[1]} \cdot x + b^{[1]} = \begin{bmatrix} 0.1 & 0.2 & 0.5 \\ 0.2 & 0.3 & 0.1 \\ 0.1 & 0.4 & 0.6 \\ 0.2 & 0.3 & 0.9 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.8 \\ 0.6 \\ 1.1 \\ 1.4 \end{bmatrix}$$

$$a^{[1]} = \text{sigmoid}(z^{[1]}) = \begin{bmatrix} 0.68 \\ 0.64 \\ 0.75 \\ 0.8 \end{bmatrix}$$

# Forward pass; Example

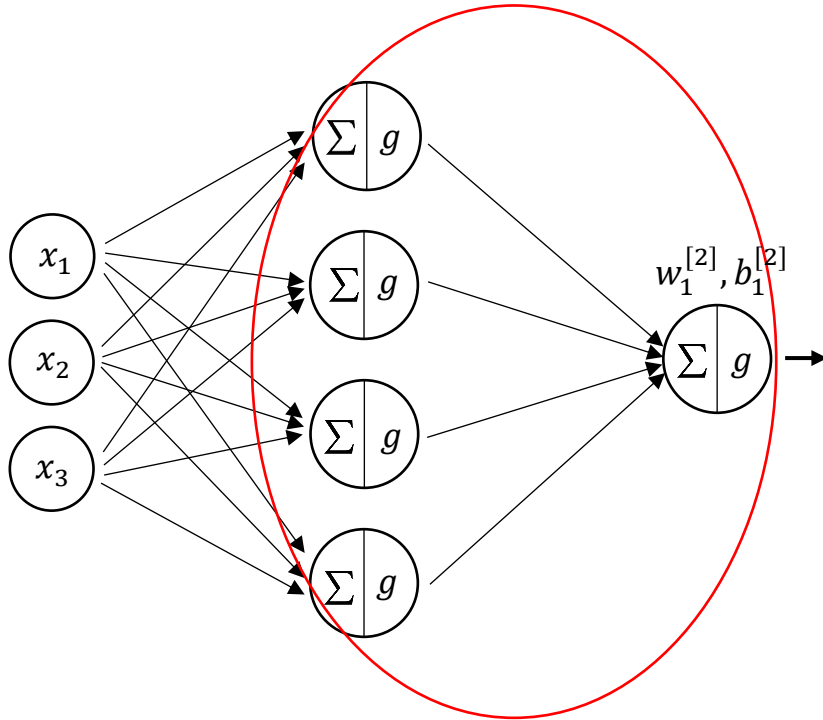


$$w^{[2]} = [0.4 \quad 1.2 \quad 4.1 \quad 0.1] \quad b^{[2]} = 1.4$$

$$g = \text{sigmoid}$$

$$\hat{y} = ?$$

# Forward pass; Example



$$w^{[2]} = [0.4 \quad 1.2 \quad 4.1 \quad 0.1] \quad b^{[2]} = 1.4$$

$$g = \text{sigmoid}$$

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

$$[0.4 \quad 1.2 \quad 4.1 \quad 0.1] \cdot \begin{bmatrix} 0.68 \\ 0.64 \\ 0.75 \\ 0.8 \end{bmatrix} + 1.4 = 5.595$$

$$a^{[2]} = \text{sigmoid}(z^{[2]}) = \mathbf{0.996}$$

$$\text{Prediction?} \quad 0.996 > 0.5 \text{ then } \hat{y} = 1$$

# Error calculation

- The **error calculation** phase of backpropagation involves **determining the difference** between the predicted output of the neural network and the desired output, which is used to quantify the error or loss.

$$\hat{y} \text{ or } P(X) \Longleftrightarrow y$$

- This error is later then **propagated backward** through the network to update the weights.
- To achieve this goal of learning weights and incorporating error, we need a **cost function**.
- ✓ What do you think cost functions should be?

Index	$y$	$P(X)$	$\hat{y}$
1	1	0.996	1
2	0	0.3	0
3	0	0.001	0
...	...	...	...
n	1	0.01	0



# Error calculation

- The specific method for calculating the error depends on the task and the type of network.
- For example, in **regression** problems, the mean squared error (**MSE**) or mean absolute error (**MAE**) may be used.

$$J(w, b) = \frac{1}{2n} \sum_{i=1}^n \left( \hat{y}^{(i)} - y^{(i)} \right)^2$$

- In **binary classification** problems, **binary cross-entropy loss** or other appropriate loss functions may be employed.

$$J(w, b) = -\frac{1}{n} \sum_{i=1}^n \left[ y^{(i)} \log \left( P_{w,b}(X^{(i)}) \right) + (1 - y^{(i)}) \log \left( 1 - P_{w,b}(X^{(i)}) \right) \right]$$

# Backward pass

- We minimize our cost function to get the best parameters. Remember that cost function is a function of parameters.

$$\text{Min } J(w, b)$$

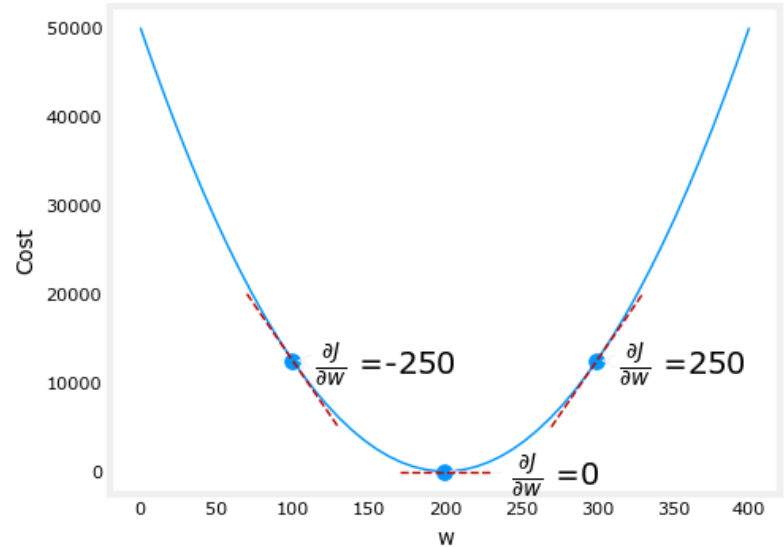
- Recall that in **Gradient Descent**:

If  $\frac{\partial J}{\partial w_i} > 0$ , then increasing  $w_i$  increases  $J$ .

If  $\frac{\partial J}{\partial w_i} < 0$ , then increasing  $w_i$  decreases  $J$ .

- So, the following update **decreases the cost function**:

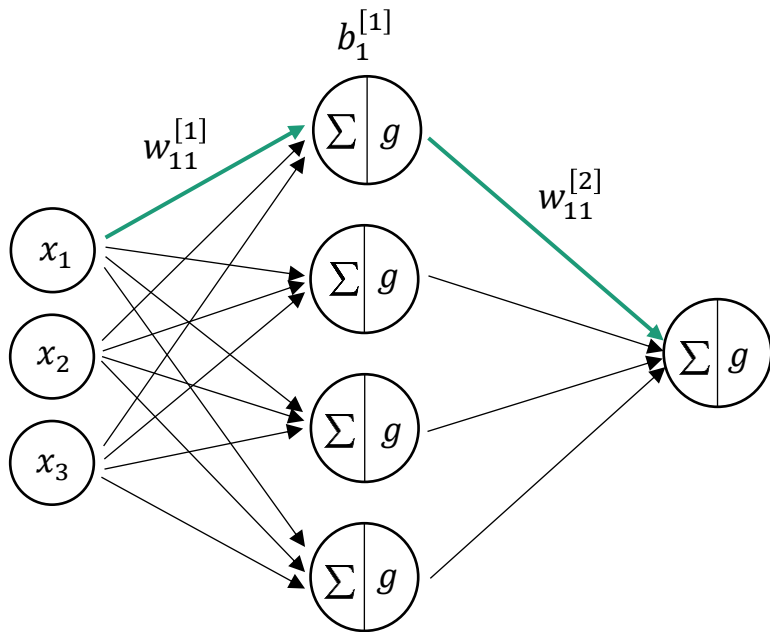
$$w_{i,\text{new}} = w_{i,\text{old}} - \alpha \frac{\partial J}{\partial w_i}$$



# Backward pass

- We must compute the gradients for each parameter.

$$w_{i,new} = w_{i,old} - \alpha \frac{\partial J}{\partial w_i}$$



$$\left\{ \begin{array}{l} w_{11,new}^{[1]} = w_{11,old}^{[1]} - \alpha \frac{\partial J}{\partial w_{11}^{[1]}} \\ b_{1,new}^{[1]} = b_{1,old}^{[1]} - \alpha \frac{\partial J}{\partial b_1^{[1]}} \\ w_{11,new}^{[2]} = w_{11,old}^{[2]} - \alpha \frac{\partial J}{\partial w_{11}^{[2]}} \end{array} \right.$$

- ✓ How to compute the gradients?

# Backward pass; Chain rule

- Example:

$$f(u) = 7u + u^2$$

$$u(x) = x^3$$

$$f(x) = f(u(x)) = 7x^3 + x^6$$

$$\frac{\partial f}{\partial x} = 21x^2 + 6x^5$$

---

$$\left. \begin{array}{l} \frac{\partial f}{\partial u} = 7 + 2u \\ \frac{\partial u}{\partial x} = 3x^2 \end{array} \right\} \begin{array}{c} \boxed{\frac{\partial f}{\partial x} = \frac{\partial f}{\partial u} \times \frac{\partial u}{\partial x}} \\ \downarrow \end{array}$$

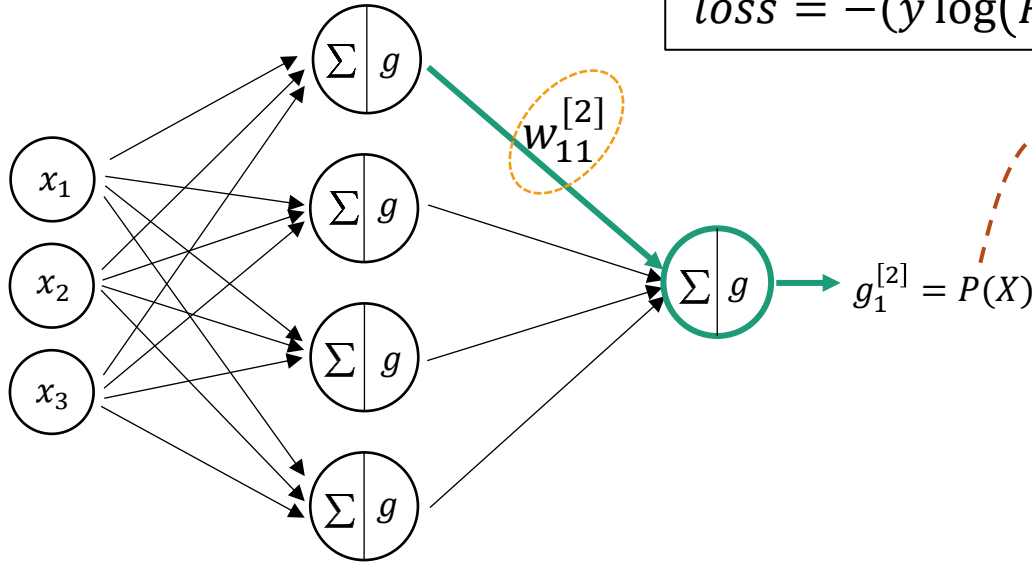
$$\frac{\partial f}{\partial x} = (7 + 2u)(3x^2) = (7 + 2x^3)(3x^2) = 21x^2 + 6x^5$$

# Backward pass; Example

- **Backpropagation** uses the **chain rule** for differentiation.

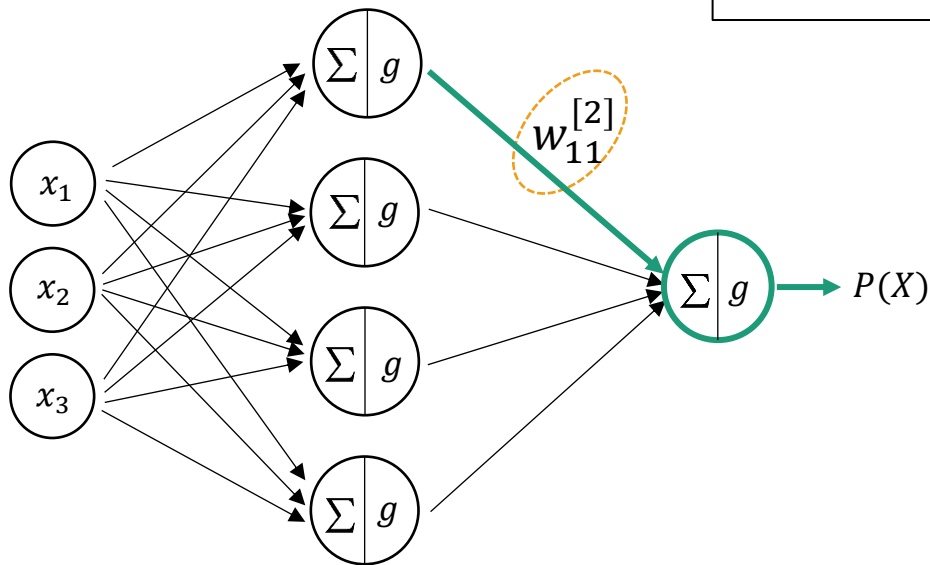
- Suppose we have one observation  $\Rightarrow n = 1$

$$loss = -(y \log(P(X)) + (1 - y) \log(1 - P(X)))$$



# Backward pass; Example

$$loss = -(y \log(P(X)) + (1 - y) \log(1 - P(X)))$$

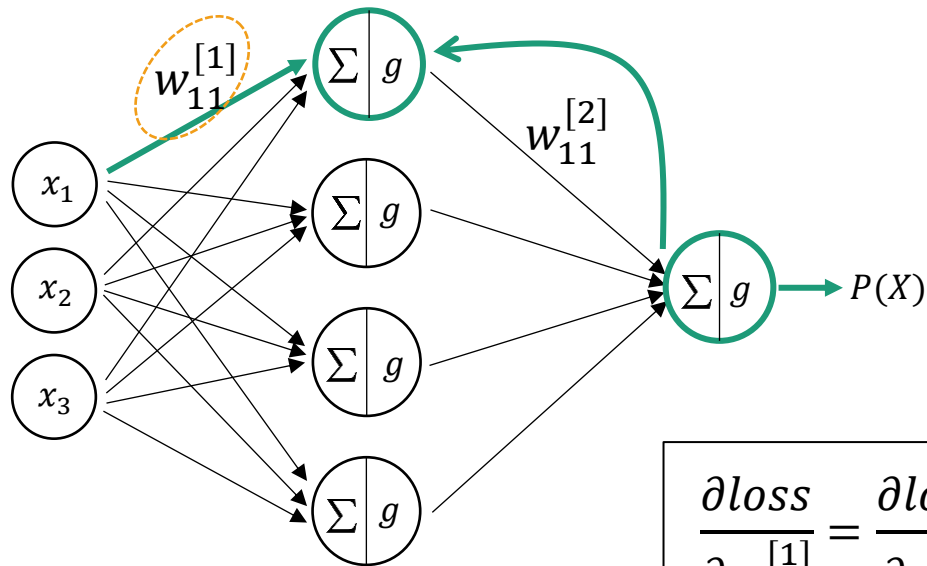


Chain rule

$$\frac{\partial loss}{\partial w_{11}^{[2]}} = \frac{\partial loss}{\partial g} \times \frac{\partial g}{\partial z_1^{[2]}} \times \frac{\partial z_1^{[2]}}{\partial w_{11}^{[2]}}$$

# Backward pass

✓ Another example:



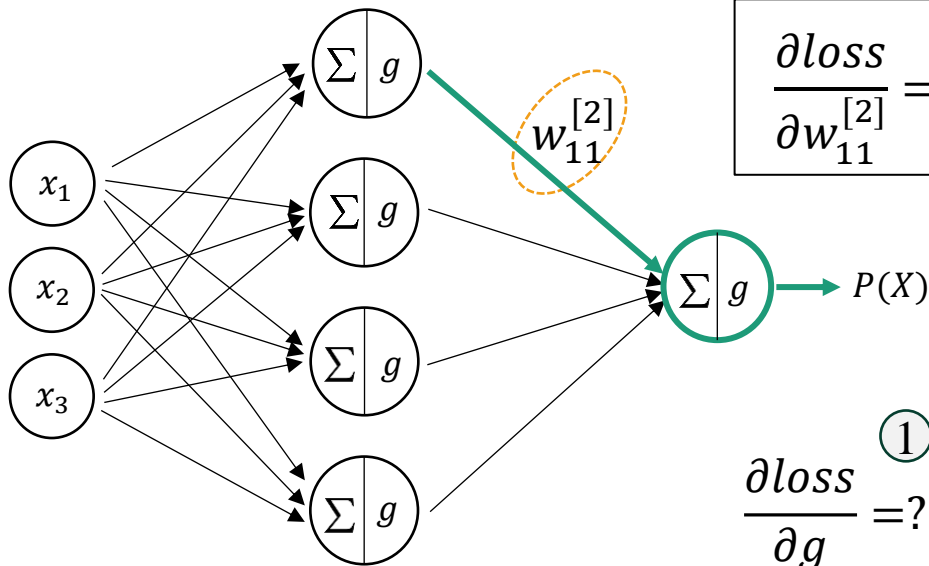
$$\text{loss} \Rightarrow g \Rightarrow z \Rightarrow w \text{ or } b$$

$$\frac{\partial \text{loss}}{\partial w_{11}^{[1]}} = \frac{\partial \text{loss}}{\partial g_1^{[2]}} \times \frac{\partial g_1^{[2]}}{\partial z_1^{[2]}} \times \frac{\partial z_1^{[2]}}{\partial g_1^{[1]}} \times \frac{\partial g_1^{[1]}}{\partial z_1^{[1]}} \times \frac{\partial z_1^{[1]}}{\partial w_{11}^{[1]}}$$

# Backward pass; Example

$$loss = -(y \log(P(X)) + (1 - y) \log(1 - P(X)))$$

$$\frac{\partial loss}{\partial w_{11}^{[2]}} = \frac{\partial loss}{\partial g} \times \frac{\partial g}{\partial z_1^{[2]}} \times \frac{\partial z_1^{[2]}}{\partial w_{11}^{[2]}}$$



$$\frac{\partial loss}{\partial g} \text{ ① } = ?$$

$$\frac{\partial g}{\partial z_1^{[2]}} \text{ ② } = ?$$

$$\frac{\partial z_1^{[2]}}{\partial w_{11}^{[2]}} \text{ ③ } = ?$$



# Differentiation; Review

$$\frac{d \log_a x}{d x} = \frac{1}{x \cdot \ln(a)}$$

$$\Rightarrow \textcircled{1} \quad \frac{\partial \text{loss}}{\partial g} = -\left(\frac{y}{g \cdot \ln(10)} + \frac{y-1}{g \cdot \ln(10)}\right) = \frac{2y-1}{g \cdot \ln(10)}$$

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

$$\textcircled{2} \quad \frac{\partial g}{\partial z_1^{[2]}} = \sigma(z_1^{[2]}) (1 - \sigma(z_1^{[2]}))$$

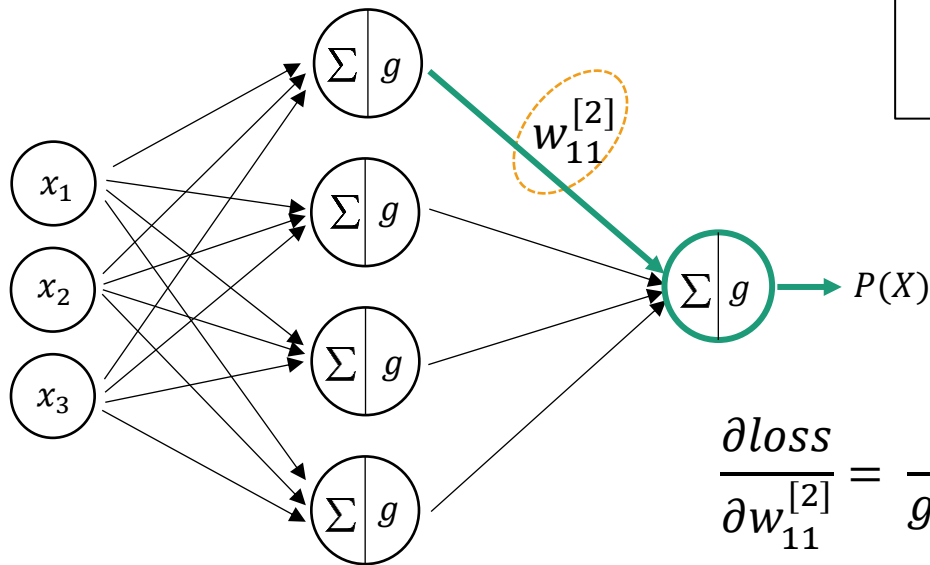
$$\frac{\partial z_j^{[k]}}{\partial w_{ij}^{[k]}} = a_i^{[k-1]}$$

$\Rightarrow$

$$\textcircled{3} \quad \frac{\partial z_1^{[2]}}{\partial w_{11}^{[2]}} = a_1^{[1]}$$

$$\frac{\partial z_j^{[k]}}{\partial b_j^{[k]}} = 1$$

# Backward pass; Example



$$\frac{\partial loss}{\partial w_{11}^{[2]}} = \frac{\partial loss}{\partial g} \times \frac{\partial g}{\partial z_1^{[2]}} \times \frac{\partial z_1^{[2]}}{\partial w_{11}^{[2]}}$$

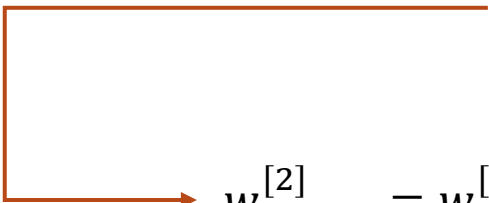
$$\frac{\partial loss}{\partial w_{11}^{[2]}} = \frac{2y - 1}{g \cdot \ln(10)} \times \sigma(z_1^{[2]}) (1 - \sigma(z_1^{[2]})) \times a_1^{[1]}$$

# Parameters updates

- We **update** all parameters (weights and biases) according to:

$$w_{i,new} = w_{i,old} - \alpha \frac{\partial J}{\partial w_i}$$

- For our example:

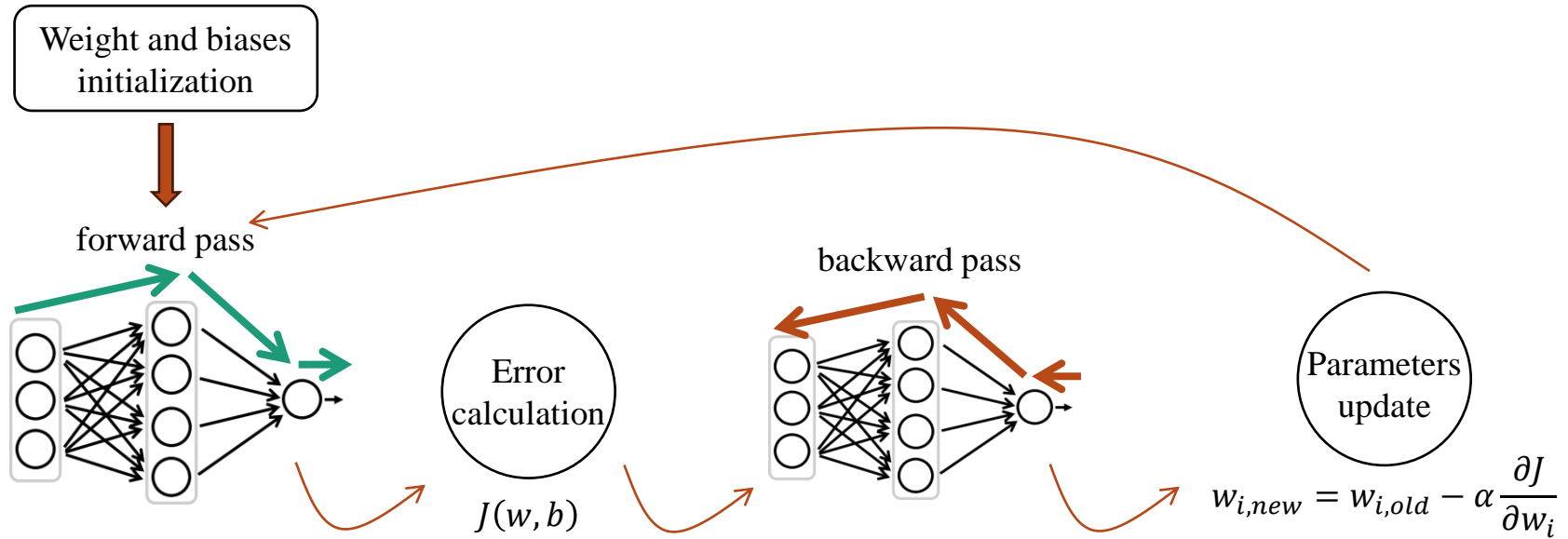

$$w_{11,new}^{[2]} = w_{11,old}^{[2]} - \alpha \frac{\partial l}{\partial w_{11}^{[2]}}$$

$$w_{11,new}^{[2]} = w_{11,old}^{[2]} - \alpha \left( \frac{2y - 1}{g \cdot \ln(10)} \times \sigma(z_1^{[2]}) (1 - \sigma(z_1^{[2]})) \times a_1^{[1]} \right)$$

- ✓ What about  $\alpha$  ?

# Backpropagation; Summary

- This process of **forward pass**, **error calculation**, **backward pass**, and **weights update** continues for multiple epochs until the network performance reaches a satisfactory level or stops improving significantly.

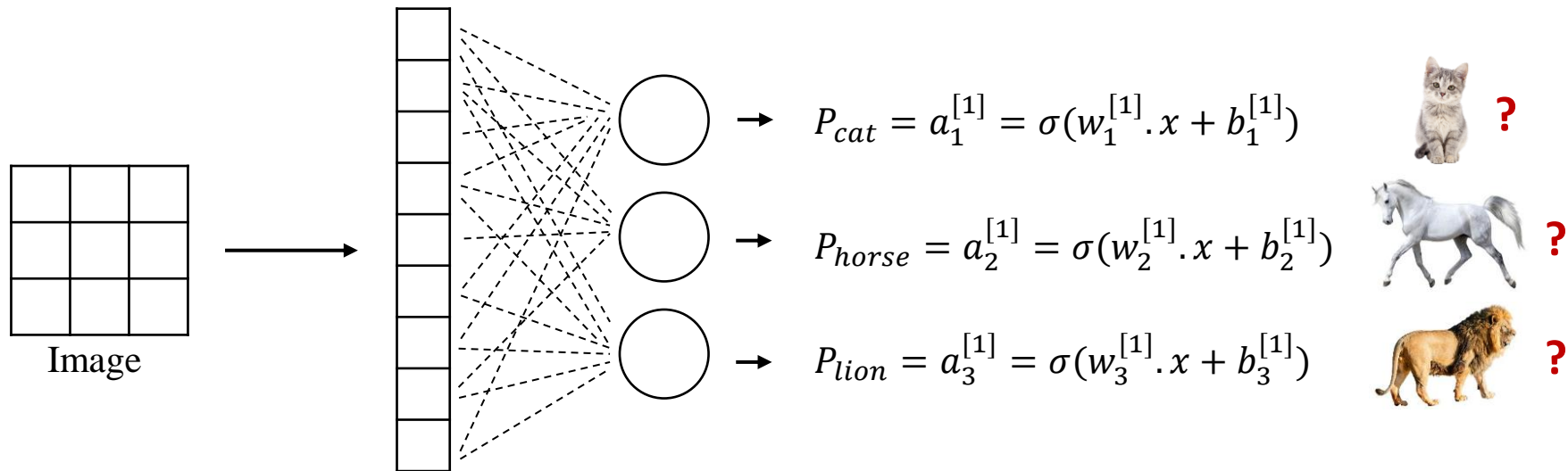


# SoftMax

---

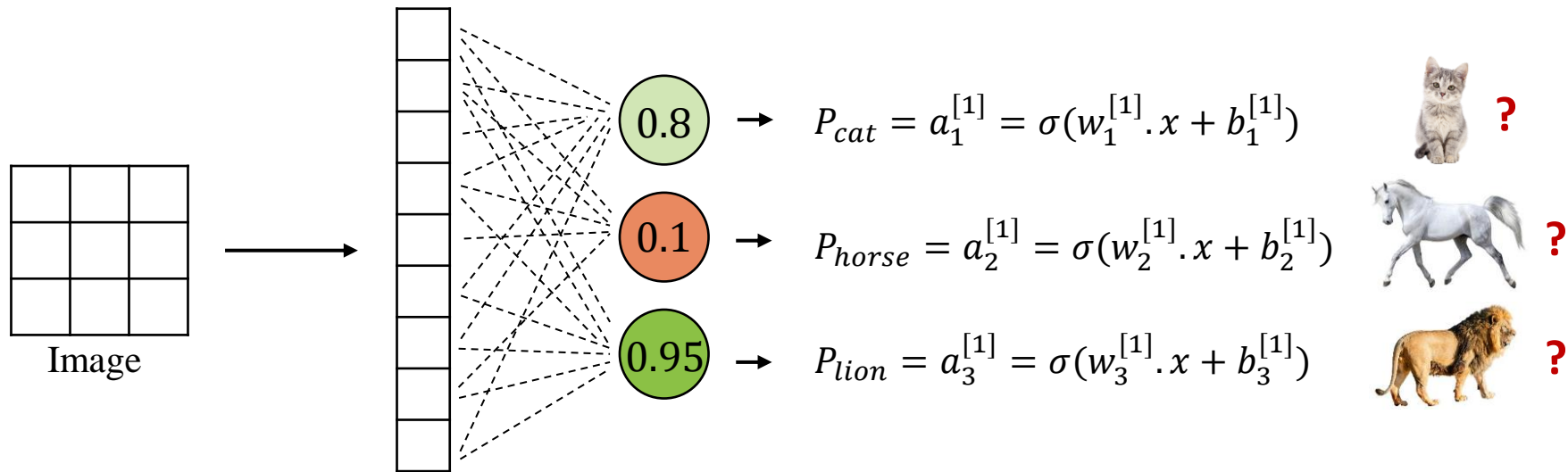
# SoftMax

- Suppose we have a **multiclass classification** problem.



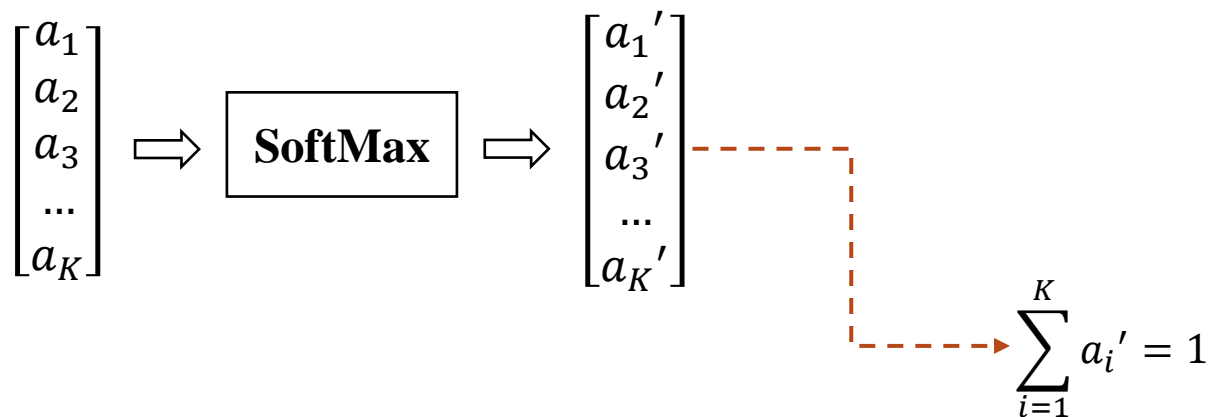
# SoftMax

✓ Considering a **threshold** of **0.5**, what label would you assign in this situation?



# SoftMax

- The **SoftMax** activation function, also known as the **normalized exponential function**, is particularly useful within the context of **multi-class classification** problems.
- The **SoftMax** function is a function that turns a vector of K real values into a vector of K real values that **sum to 1**.

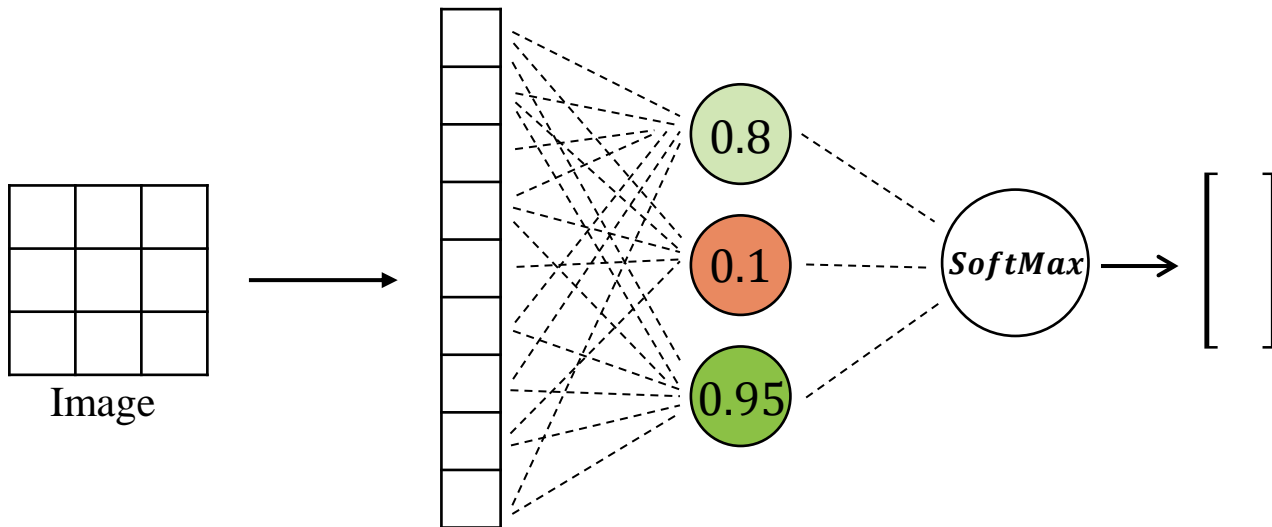




# SoftMax

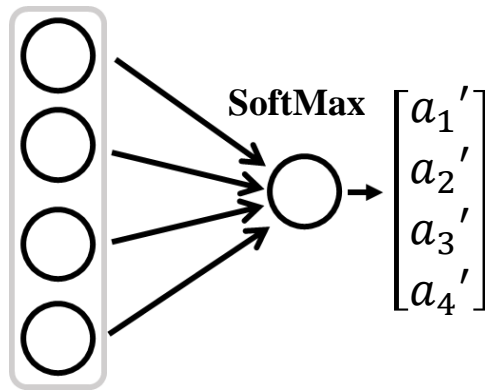
- For input vector  $z$  with elements  $z_1, z_2, \dots, z_k$  the **SoftMax** function is defined as:

$$g(z_i) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$



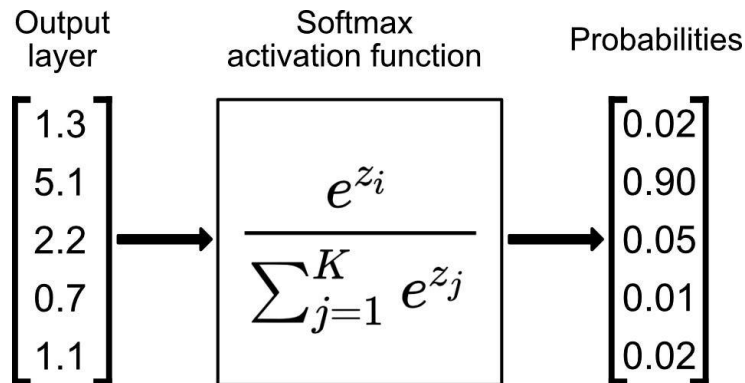
# SoftMax; Architecture

- For the **SoftMax** layer to function correctly, we need **a neuron for each class** in the previous layer.
- Consequently, the SoftMax output will also be a vector of this size.
- For example, in a four-class classification:



# SoftMax

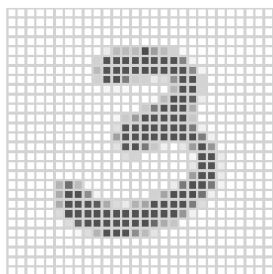
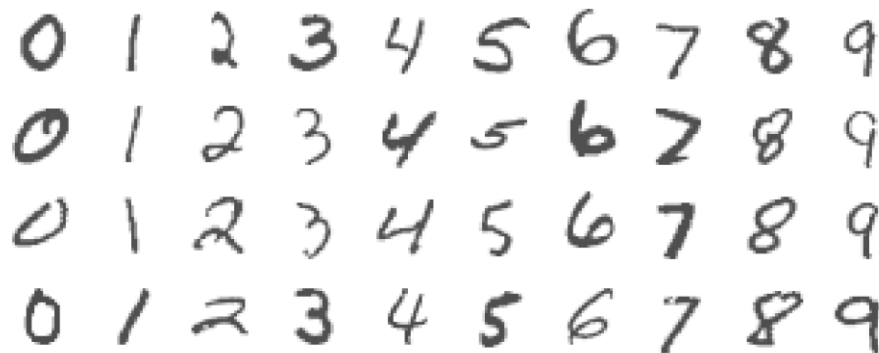
- The input values can be positive, negative, zero, or greater than one, but the **SoftMax** transforms them into values between 0 and 1, so that they can be interpreted as probabilities.



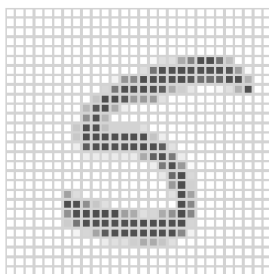
- There's no necessity to use another activation function like **sigmoid** before **SoftMax**. In fact, Since **SoftMax** already handles the conversion to probabilities, using sigmoid beforehand is redundant in the context of multi-class classification.

# Multi-class Classification; example

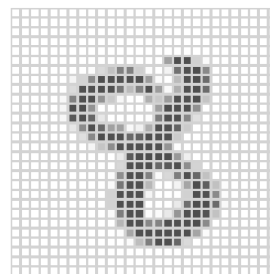
- [MNIST](#) is digit recognition dataset.
  - Handwritten digits
  - $28 \times 28$  grayscale images
  - Grayscale values:  $[0, 255]$
  - Features are the 784 pixel
  - 60K train, 10K test images
  - Labels are the digit class 0 to 9



3



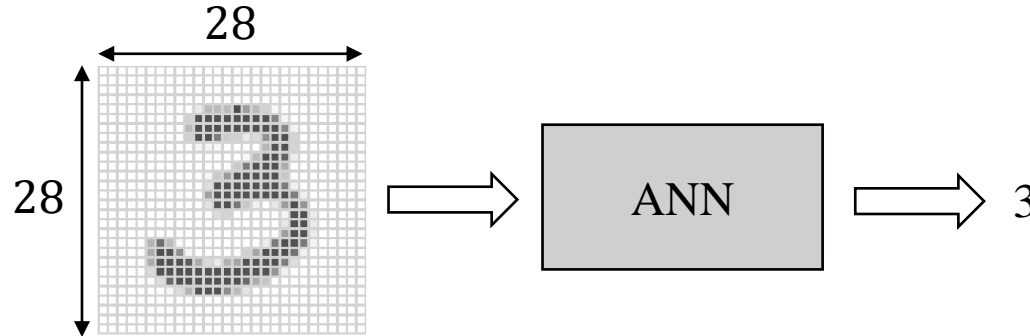
5



8

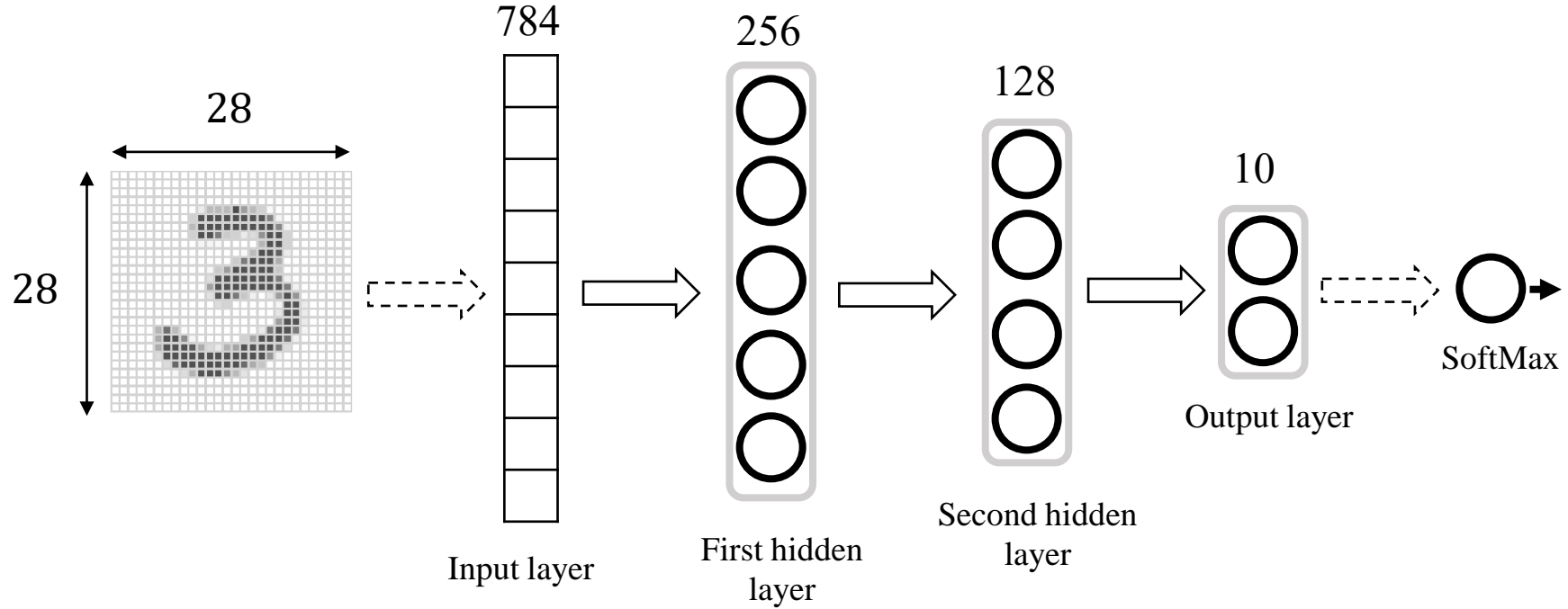
# Multi-class Classification; example

- Our goal is to build a classifier to predict the image class.
- We build a **two-layer network** with 256 units at first layer, 128 units at second layer, and 10 units at output layer.



- Along with biases, there are **235,146 parameters**.
- ✓ Comparing number of parameters with previous simpler models, what can you conclude?

# Multi-class Classification; example



# Multi-class Classification; example

- We fit the model by minimizing **cross-entropy**.

$$J(w, b) = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K y_j^{(i)} \cdot \log(P_{w,b}(X^{(i)}))$$

- Model performance: **error rate** < **0.5%**
- For comparison, human error rate is reported to be around **0.2%**, or 20 of the 10K test images.
- Some of hard examples in MNIST: \_\_\_\_\_

