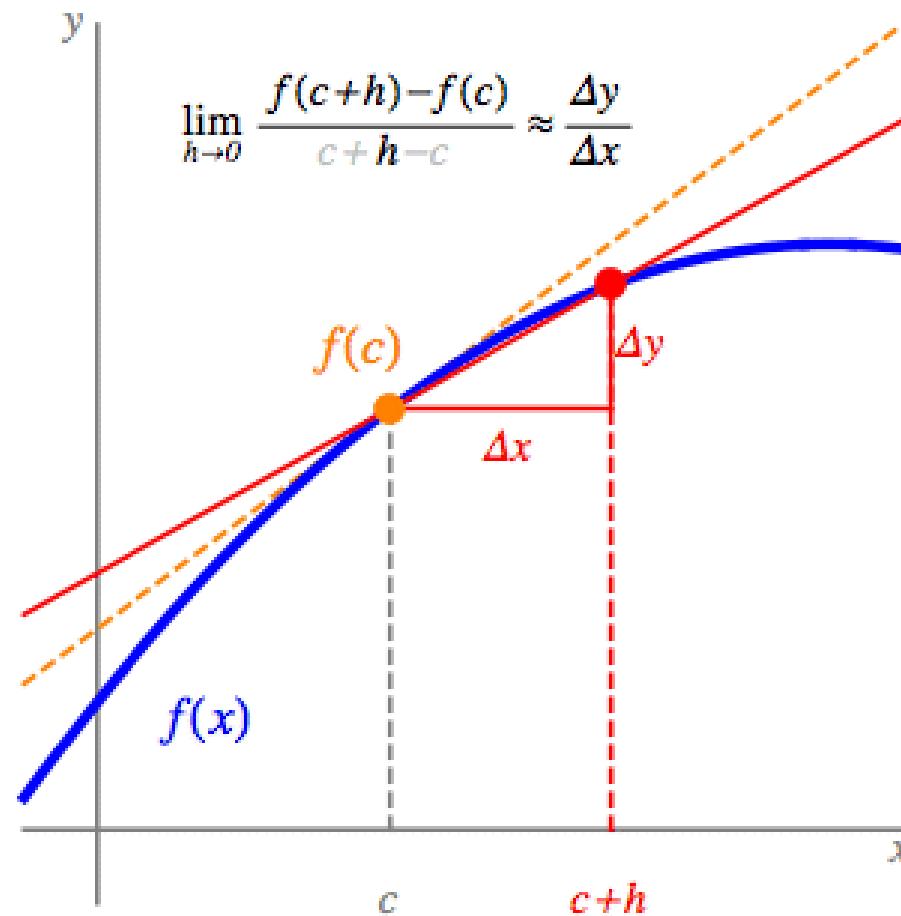


## *Limit Definition of the Derivative $f'(c)$*



[https://github.com/LambdaSchool/DS-Unit-4-Sprint-2-Neural-Networks/blob/main/module2-Train/LS\\_DS\\_422\\_Train\\_Lecture.ipynb](https://github.com/LambdaSchool/DS-Unit-4-Sprint-2-Neural-Networks/blob/main/module2-Train/LS_DS_422_Train_Lecture.ipynb)

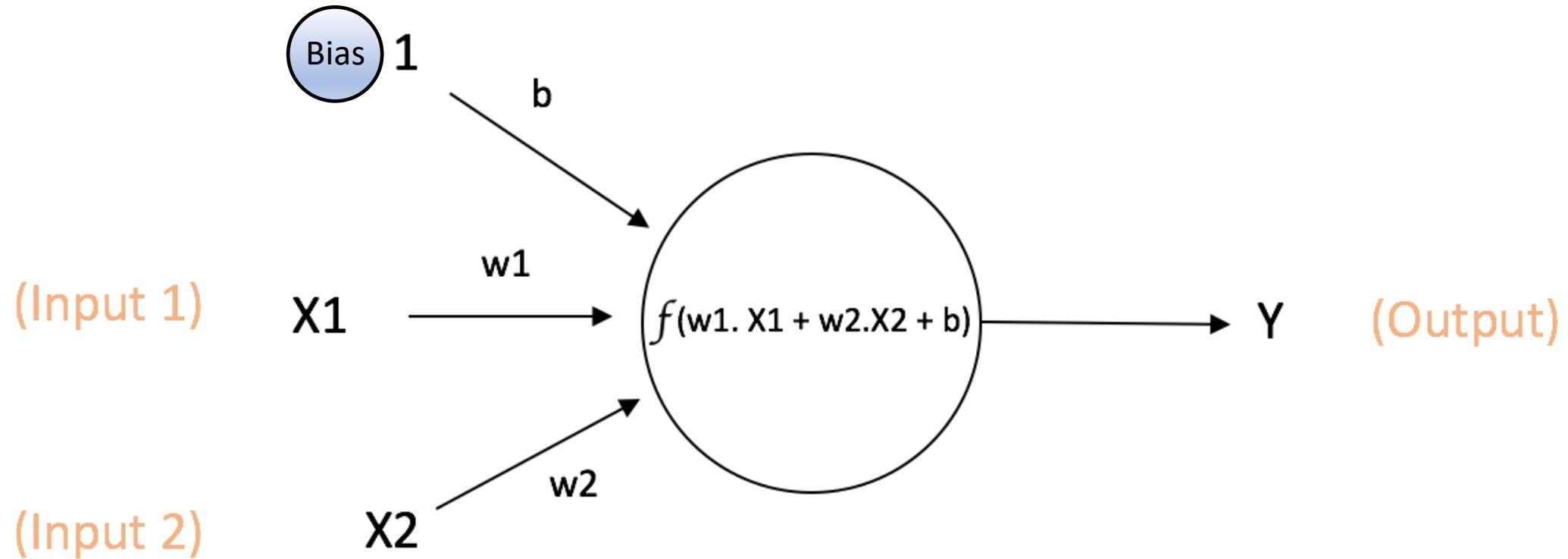
# Review of Matrix Multiplication (linear algebra)

Dimension:  
 $m \times n$        $n \times t$        $m \times t$

$$= \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \cdot \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}$$
$$= \begin{bmatrix} 1(5) + 2(7) & 1(6) + 2(8) \\ 3(5) + 4(7) & 3(6) + 4(8) \end{bmatrix}$$
$$= \begin{bmatrix} 19 & 22 \\ 43 & 50 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 1 \\ 4 & 3 \\ 5 & -1 \end{bmatrix}_{3 \times 2} \begin{bmatrix} 2 & -1 & 0 \\ 3 & 4 & 1 \end{bmatrix}_{2 \times 3} = \begin{bmatrix} 0+3 & 0+4 & 0+1 \\ 8+9 & -4+12 & 0+3 \\ 10-3 & -5-4 & 0-1 \end{bmatrix} = \begin{bmatrix} 3 & 4 & 1 \\ 17 & 8 & 3 \\ 7 & -9 & -1 \end{bmatrix}$$

A Good [References](#)



$$\text{Output of neuron} = Y = f(w_1 \cdot X_1 + w_2 \cdot X_2 + b)$$

## Network Weight Initialization [¶](#)

how we initialize our model weights can determine the difference between Gradient Descent converging towards a local minimum or a global minimum!

[Keras has documentation on all these options](#)

```
init_mode = ['uniform', 'lecun_uniform', 'normal', 'zero', 'glorot_normal',
'glorot_uniform', 'he_normal', 'he_uniform']
```

If you wish to dive into a deep analysis of varying weight initializers and their affect on model performance [check this article.](#)(optional).

**Take Away:** Include a few different weight initializers in your gridsearch.

# Weight Initialization

Equation 11-1. Xavier initialization (when using the logistic activation function)

Normal distribution with mean 0 and standard deviation  $\sigma = \sqrt{\frac{2}{n_{\text{inputs}} + n_{\text{outputs}}}}$

Or a uniform distribution between  $-r$  and  $+r$ , with  $r = \sqrt{\frac{6}{n_{\text{inputs}} + n_{\text{outputs}}}}$

Table 11-1. Initialization parameters for each type of activation function

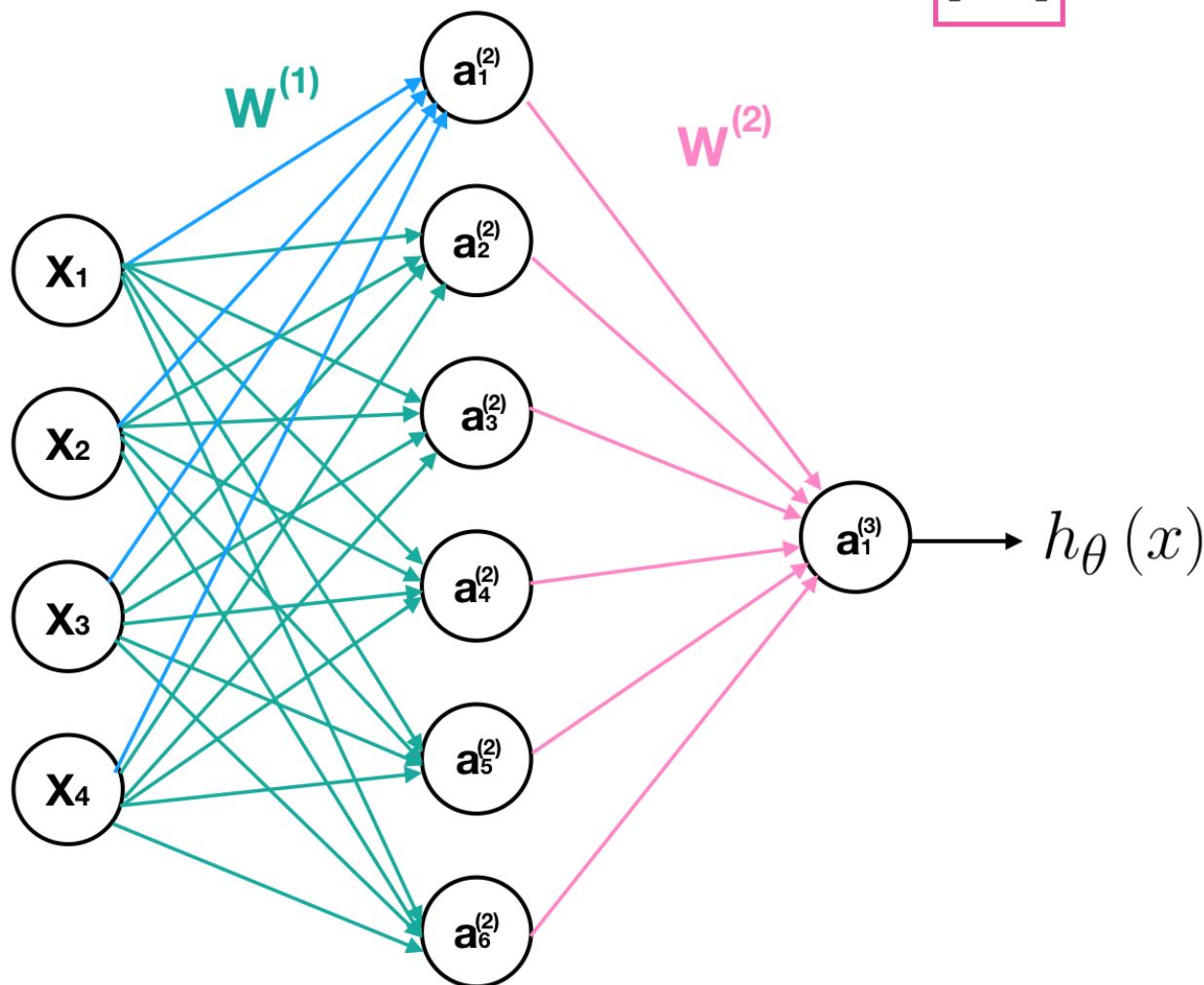
| Activation function     | Uniform distribution $[-r, r]$  | Normal distribution  |
|-------------------------|---|--|
| Logistic                | $r = \sqrt{\frac{6}{n_{\text{inputs}} + n_{\text{outputs}}}}$         | $\sigma = \sqrt{\frac{2}{n_{\text{inputs}} + n_{\text{outputs}}}}$         |
| Hyperbolic tangent      | $r = 4\sqrt{\frac{6}{n_{\text{inputs}} + n_{\text{outputs}}}}$        | $\sigma = 4\sqrt{\frac{2}{n_{\text{inputs}} + n_{\text{outputs}}}}$        |
| ReLU (and its variants) | $r = \sqrt{2}\sqrt{\frac{6}{n_{\text{inputs}} + n_{\text{outputs}}}}$ | $\sigma = \sqrt{2}\sqrt{\frac{2}{n_{\text{inputs}} + n_{\text{outputs}}}}$ |

"This initialization strategy is often called *Xavier initialization* (after the author's first name), or sometimes *Glorot initialization*."  
*Hands-On Machine Learning with Scikit-Learn and TensorFlow* Aurélien Géron 2017

# Forward Propagation

$$W^{(1)} = \begin{bmatrix} \theta_{11}^{(1)} & \theta_{21}^{(1)} & \theta_{31}^{(1)} & \theta_{41}^{(1)} & \theta_{51}^{(1)} & \theta_{61}^{(1)} \\ \theta_{12}^{(1)} & \theta_{22}^{(1)} & \theta_{32}^{(1)} & \theta_{42}^{(1)} & \theta_{52}^{(1)} & \theta_{62}^{(1)} \\ \theta_{13}^{(1)} & \theta_{23}^{(1)} & \theta_{33}^{(1)} & \theta_{43}^{(1)} & \theta_{53}^{(1)} & \theta_{63}^{(1)} \\ \theta_{14}^{(1)} & \theta_{24}^{(1)} & \theta_{34}^{(1)} & \theta_{44}^{(1)} & \theta_{54}^{(1)} & \theta_{64}^{(1)} \end{bmatrix}$$

$$W^{(2)} = \begin{bmatrix} \theta_{11}^{(2)} \\ \theta_{12}^{(2)} \\ \theta_{13}^{(2)} \\ \theta_{14}^{(2)} \\ \theta_{15}^{(2)} \\ \theta_{16}^{(2)} \end{bmatrix}$$



# Forward Propagation

$$W^T X = \begin{bmatrix} \theta_{11}^{(1)} & \theta_{12}^{(1)} & \theta_{13}^{(1)} & \theta_{14}^{(1)} \\ \theta_{21}^{(1)} & \theta_{22}^{(1)} & \theta_{23}^{(1)} & \theta_{24}^{(1)} \\ \theta_{31}^{(1)} & \theta_{32}^{(1)} & \theta_{33}^{(1)} & \theta_{34}^{(1)} \\ \theta_{41}^{(1)} & \theta_{42}^{(1)} & \theta_{43}^{(1)} & \theta_{44}^{(1)} \\ \theta_{51}^{(1)} & \theta_{52}^{(1)} & \theta_{53}^{(1)} & \theta_{54}^{(1)} \\ \theta_{61}^{(1)} & \theta_{62}^{(1)} & \theta_{63}^{(1)} & \theta_{64}^{(1)} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{bmatrix} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \\ z_5^{(1)} \\ z_6^{(1)} \end{bmatrix} = Z^{(2)}$$

$$z_1^{(2)} = \theta_{11}^{(1)}x_1 + \theta_{12}^{(1)}x_2 + \theta_{13}^{(1)}x_3 + \theta_{14}^{(1)}x_4$$

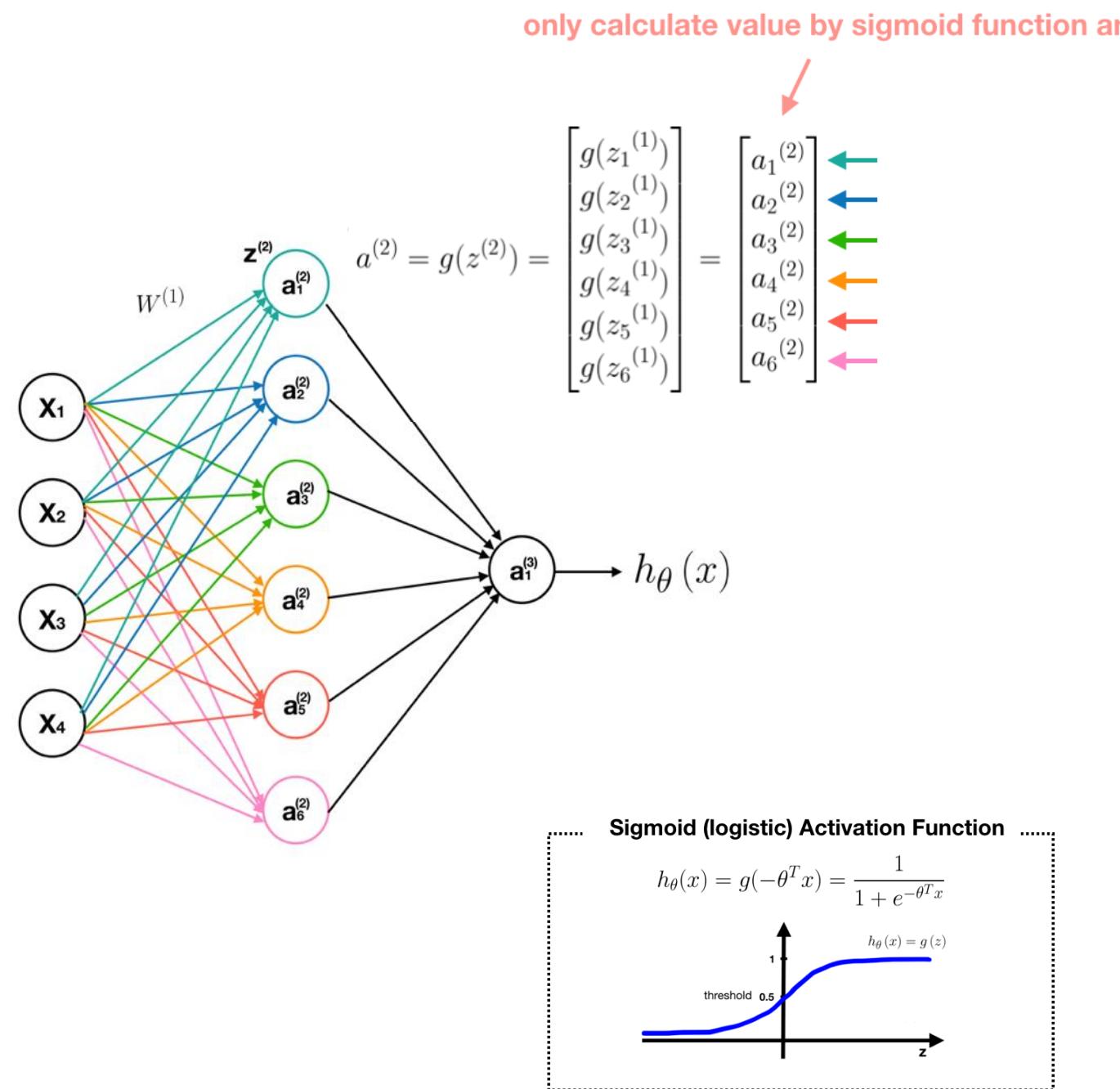
$$z_2^{(2)} = \theta_{21}^{(1)}x_1 + \theta_{22}^{(1)}x_2 + \theta_{23}^{(1)}x_3 + \theta_{24}^{(1)}x_4$$

$$z_3^{(2)} = \theta_{31}^{(1)}x_1 + \theta_{32}^{(1)}x_2 + \theta_{33}^{(1)}x_3 + \theta_{34}^{(1)}x_4$$

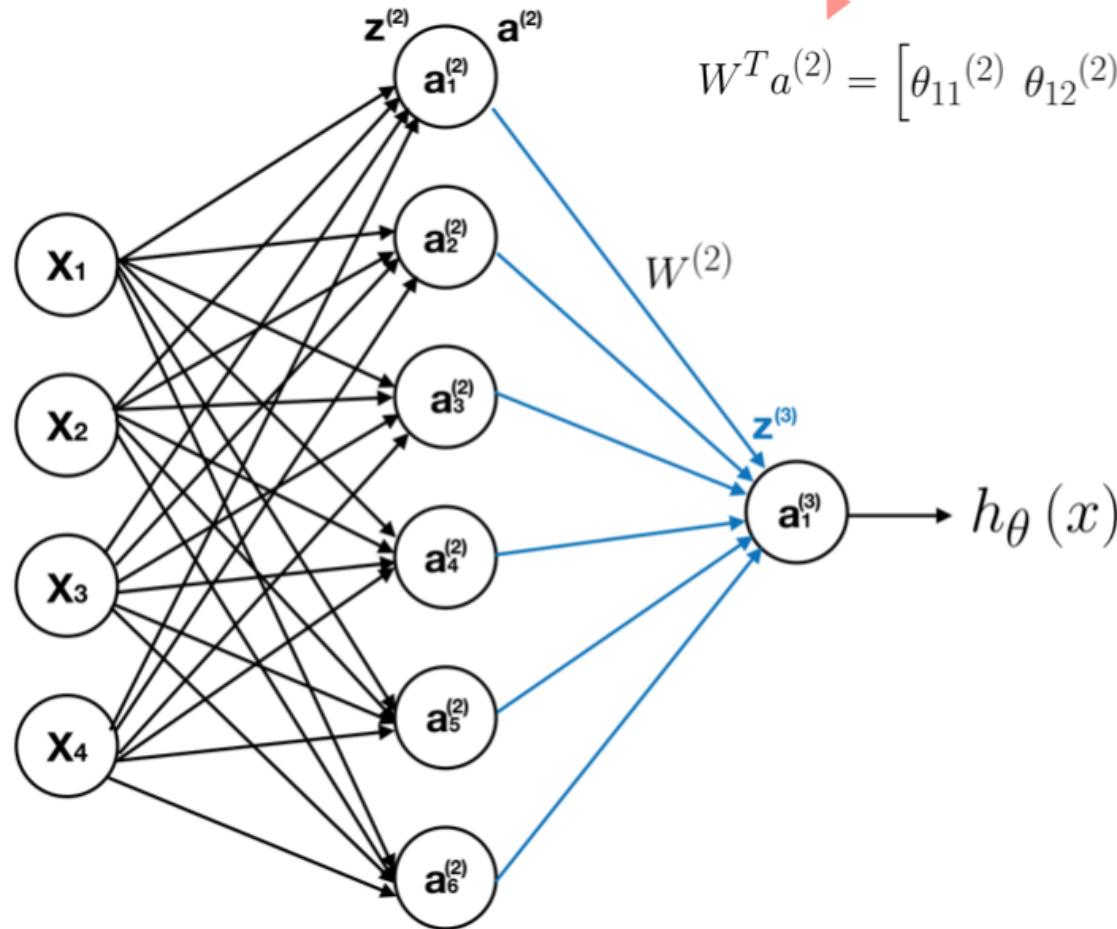
$$z_4^{(2)} = \theta_{41}^{(1)}x_1 + \theta_{42}^{(1)}x_2 + \theta_{43}^{(1)}x_3 + \theta_{44}^{(1)}x_4$$

$$z_5^{(2)} = \theta_{51}^{(1)}x_1 + \theta_{52}^{(1)}x_2 + \theta_{53}^{(1)}x_3 + \theta_{54}^{(1)}x_4$$

$$z_6^{(2)} = \theta_{61}^{(1)}x_1 + \theta_{62}^{(1)}x_2 + \theta_{63}^{(1)}x_3 + \theta_{64}^{(1)}x_4$$



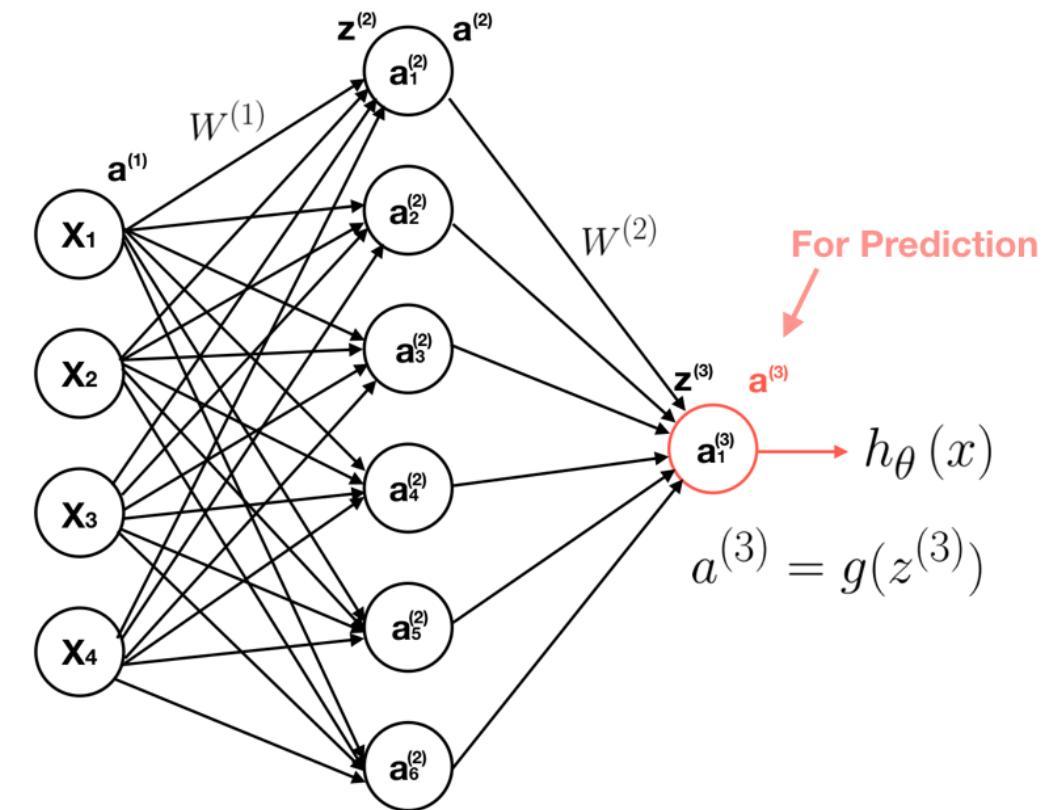
# Forward Propagation



Passing Value

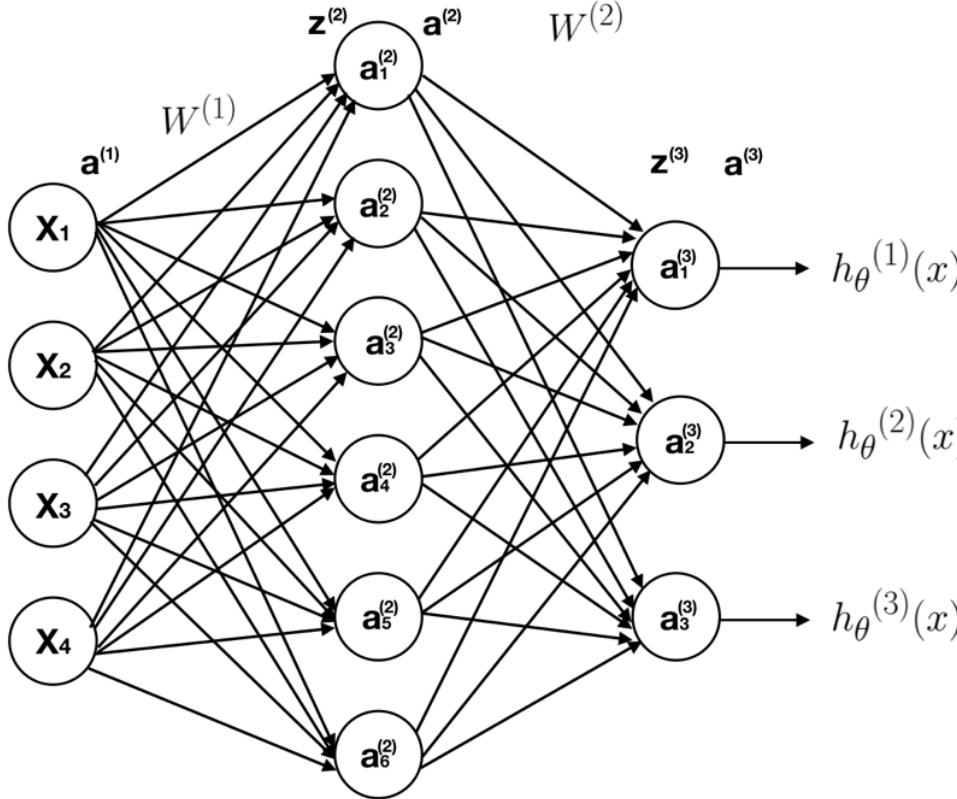
$$W^T a^{(2)} = \left[ \theta_{11}^{(2)} \ \theta_{12}^{(2)} \ \theta_{13}^{(2)} \ \theta_{14}^{(2)} \ \theta_{15}^{(2)} \ \theta_{16}^{(2)} \right] = z^{(3)}$$

$$\begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \\ a_4^{(2)} \\ a_5^{(2)} \\ a_6^{(2)} \end{bmatrix} = z^{(3)}$$



For Prediction

# Forward Propagation



## Forward Propagation

$$W^{(1)T} X = z^{(2)}$$

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

$$W^{(2)T} a^{(2)} = z^{(3)}$$

$$a^{(3)} = g(z^{(3)})$$

For Prediction

$$a^{(3)} = \begin{bmatrix} a_1^{(3)} \\ a_2^{(3)} \\ a_3^{(3)} \end{bmatrix} \rightarrow \begin{array}{l} P(x \in Class1) \\ P(x \in Class2) \\ P(x \in Class3) \end{array}$$

For Prediction

Multi-Classification Problem

# Forward Propagation: Gradient Descent

$$J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\Theta(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_\Theta(x^{(i)}))_k) \right] + \boxed{\frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2}$$

↓  
data number  
↑  
label number

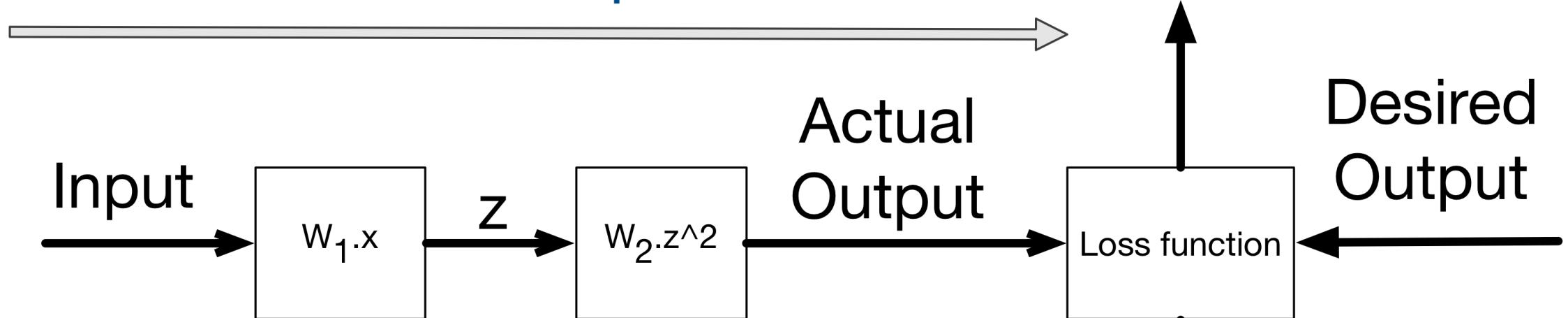
*Repeat until converge {*

$$\Theta_{ij}^{(l)} = \Theta_{ij}^{(l)} - \boxed{\alpha} \frac{\partial J(\Theta)}{\partial \Theta_{ij}^{(l)}}$$

*}*

**learning rate**

## Forward pass



$dW_1$

$dW_2$

$w_1$

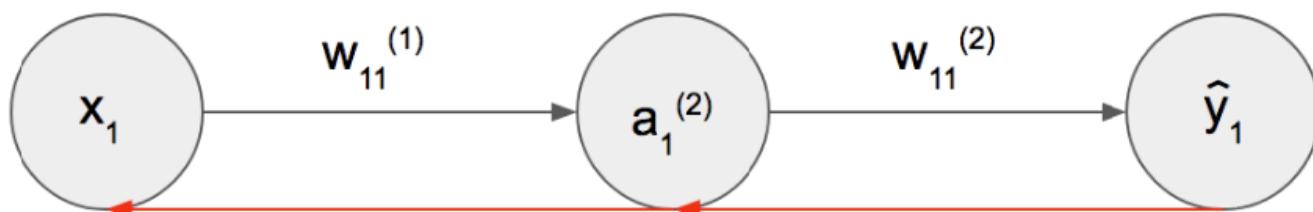
$2 \cdot w_2 \cdot z$

Derivative  
of loss

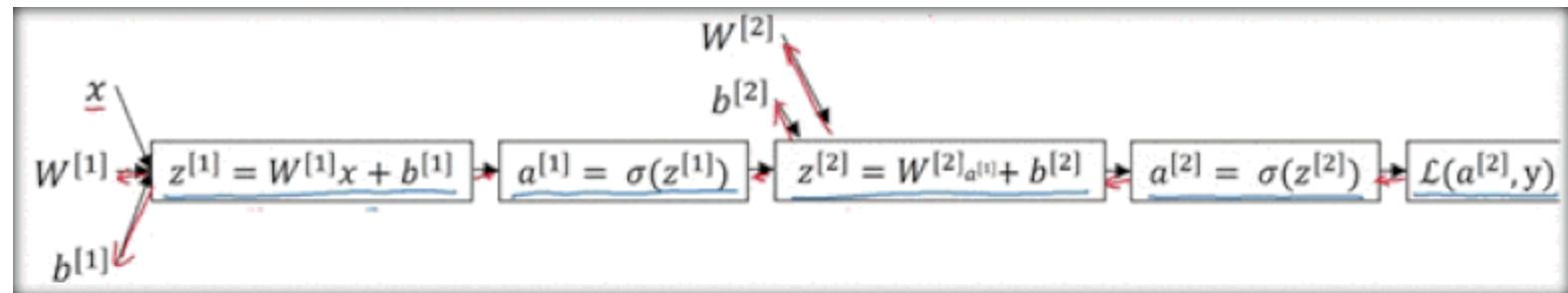
$dz$

$dE$

Back-propagate error



$$\frac{\partial E}{\partial w_{11}^{(1)}} = \frac{\partial E}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial w_{11}^{(1)}} \quad \frac{\partial E}{\partial w_{11}^{(2)}} = \frac{\partial E}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial w_{11}^{(2)}}$$



$$\frac{\partial L}{\partial a}$$

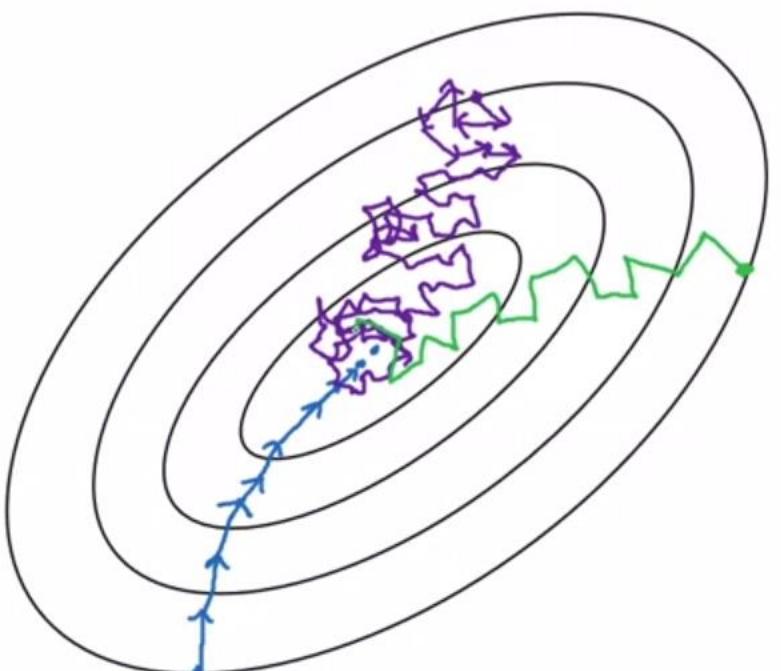
All

# Choosing your mini-batch size

→ If mini-batch size = m : Batch gehört versch.

→ If min-batch size = 1 : Stochastic gradient descent. Every example is its own  $(X^{(i)}, Y^{(i)}) = (x^{(i)}, y^{(i)}) \dots (x^{(n)}, y^{(n)})$  min-batch.

In practice: Some in-between l and m



Stochastic  
gradient  
descent

Use speakers  
from ventilation

In-between  
(minihatch size  
not too big/small)

## Faster learning.

- Vectorization.  
( $n \geq 1000$ )
  - Make passes without  
processing entire query set.

# Batch size

$$(X^{\{1\}}, Y^{\{1\}}) = (X, Y)$$

Batch  
gradient descent  
(mini-batch size =  $n$ )

↓  
Too long  
per iteration

# Choosing your mini-batch size

If small training set : Use batch gradient descent.  
 $(m \leq 2000)$

Typical mini-batch sizes:

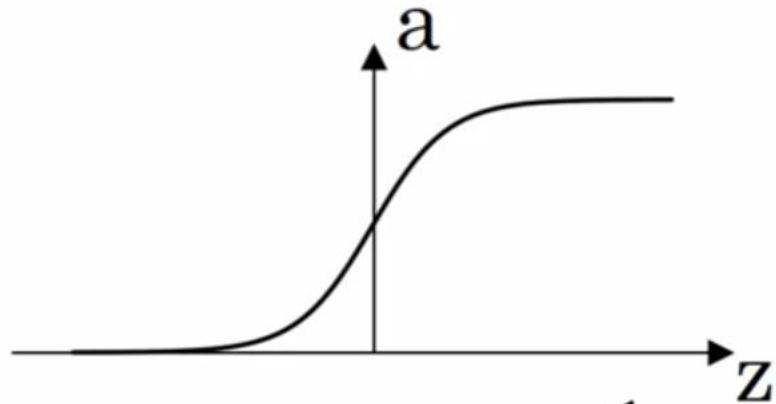
$$\underbrace{64, 128, 256, 512}_{2^6, 2^7, 2^8, 2^9} \quad \frac{1024}{2^{10}}$$

Make sure mini-batch fits in CPU/GPU memory.

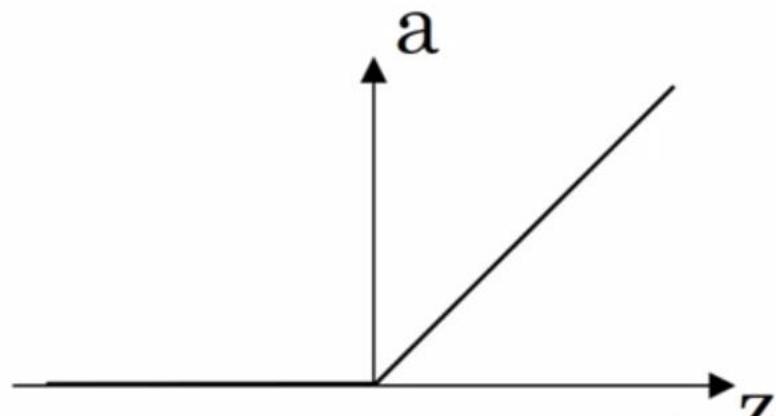
$$X^{\{t\}}, Y^{\{t\}}$$

# Pros and cons of activation functions

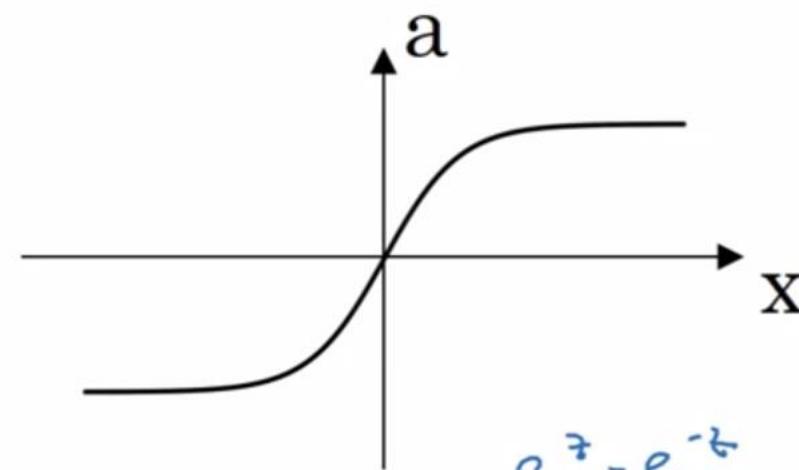
Activation Function



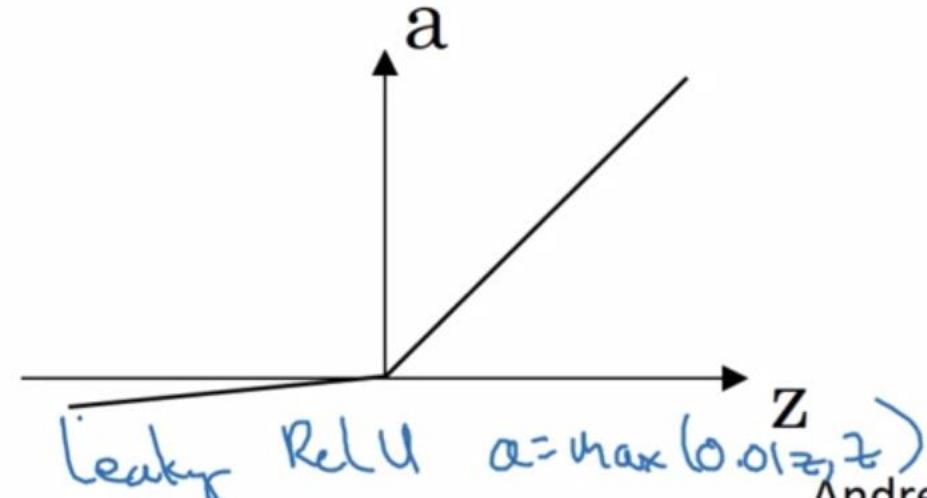
$$\text{sigmoid: } a = \frac{1}{1 + e^{-z}}$$



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



$$\tanh: \quad a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

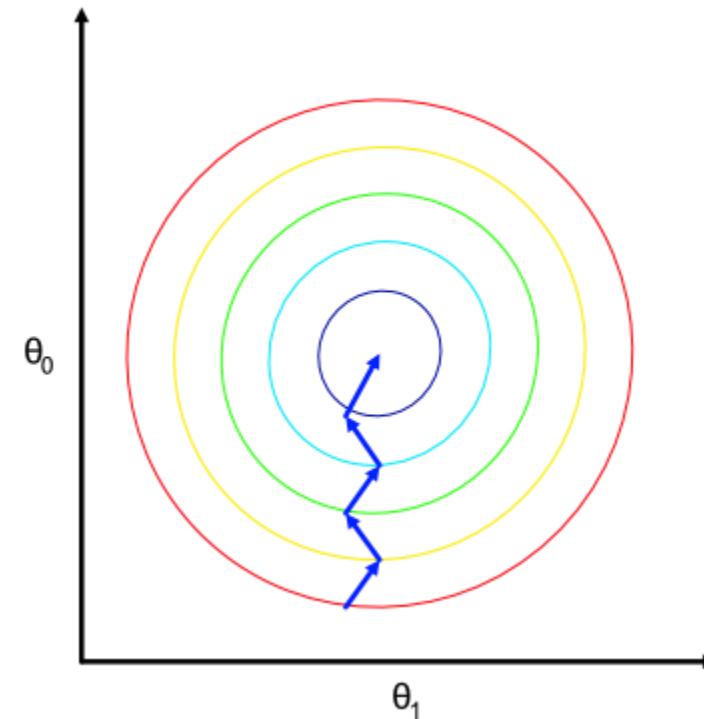
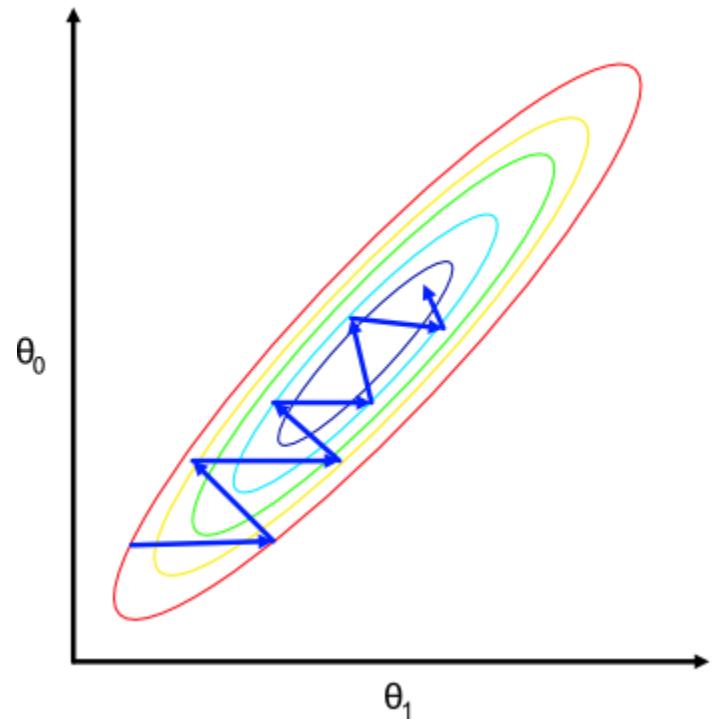


$$\text{Leaky ReLU} \quad a = \max(0, 0.01z)$$

Andrew Ng

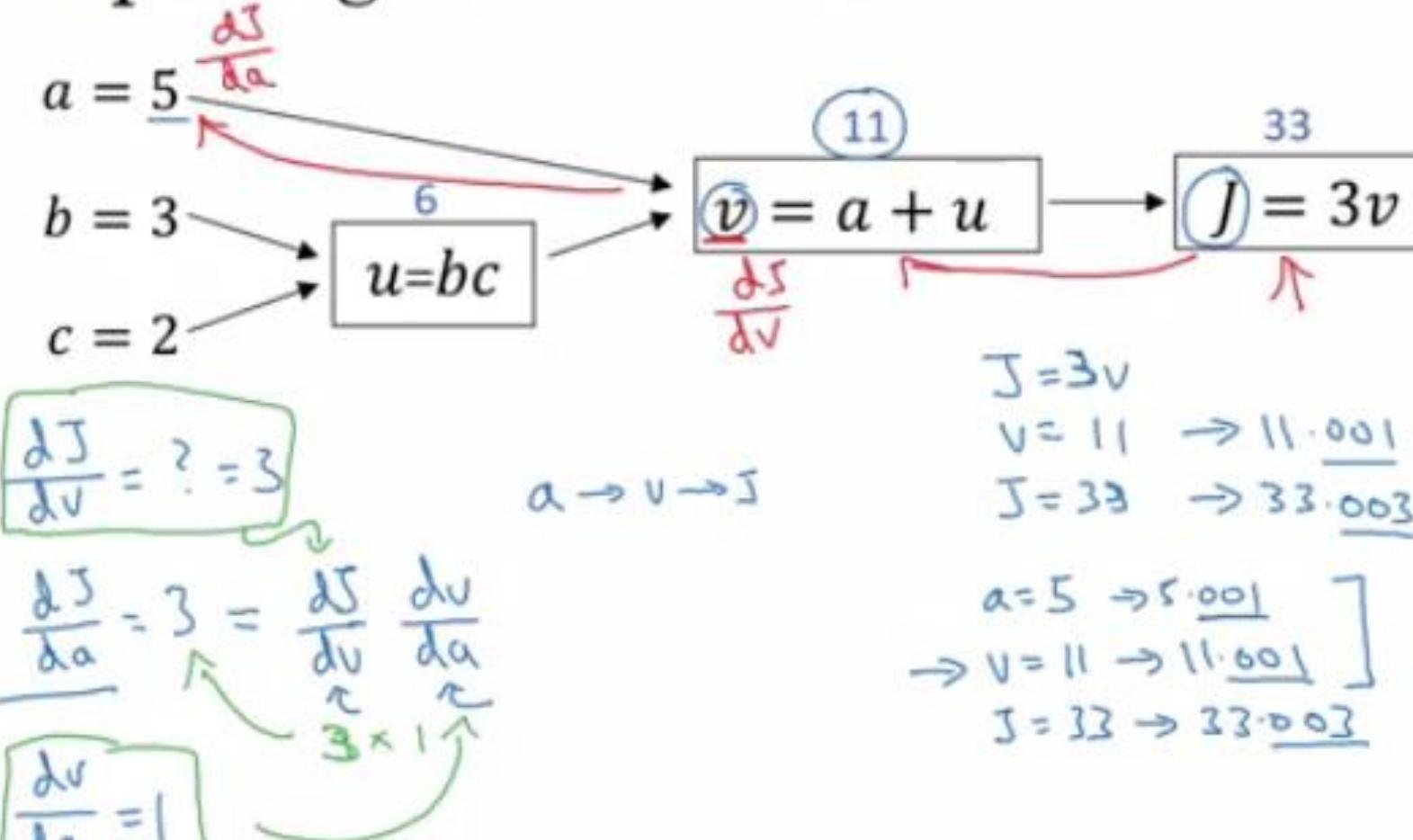
```
model.compile(loss='binary_crossentropy',
              optimizer='adam',
              metrics=['accuracy'])
model.fit(X, y, epochs=5)
```

# (Inputs) Normalization



# Appendix

# Computing derivatives



$f(a) = 3a$   
 $\frac{df(w)}{da} = \frac{df}{da} = 3$   
 $J = 3v$   
 $\frac{dJ}{dv} = 3$

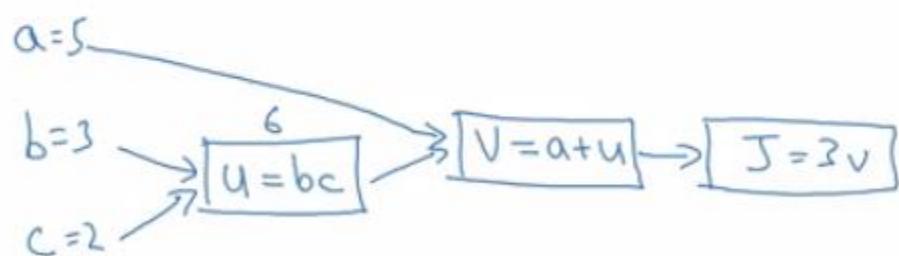
## Computation Graph

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

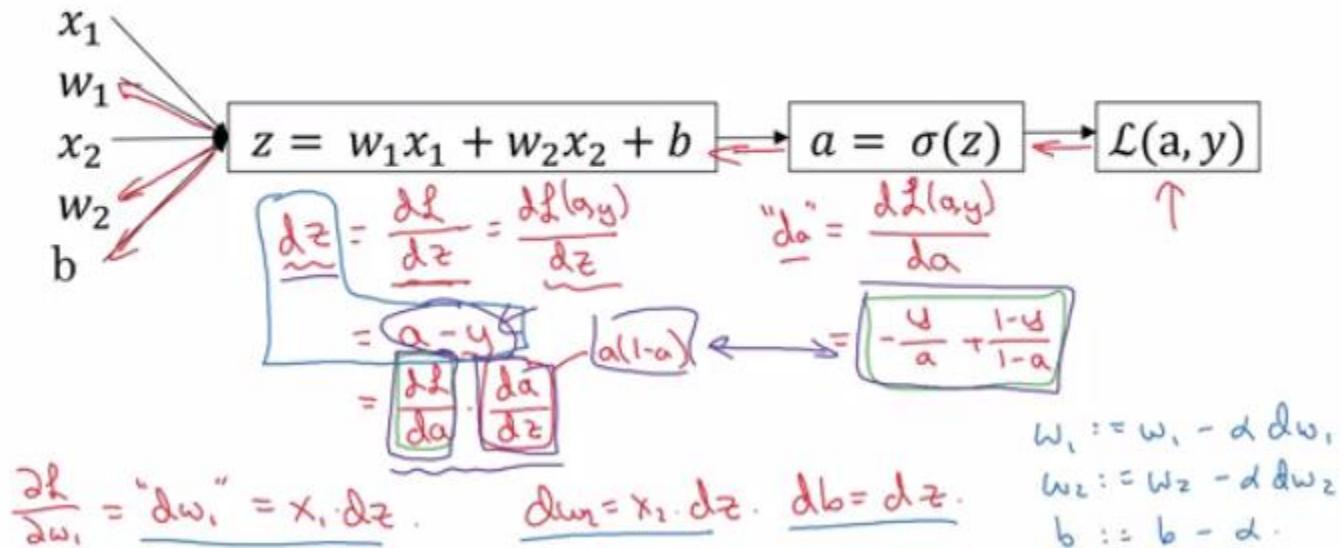
$$u = bc$$

$$v = a + u$$

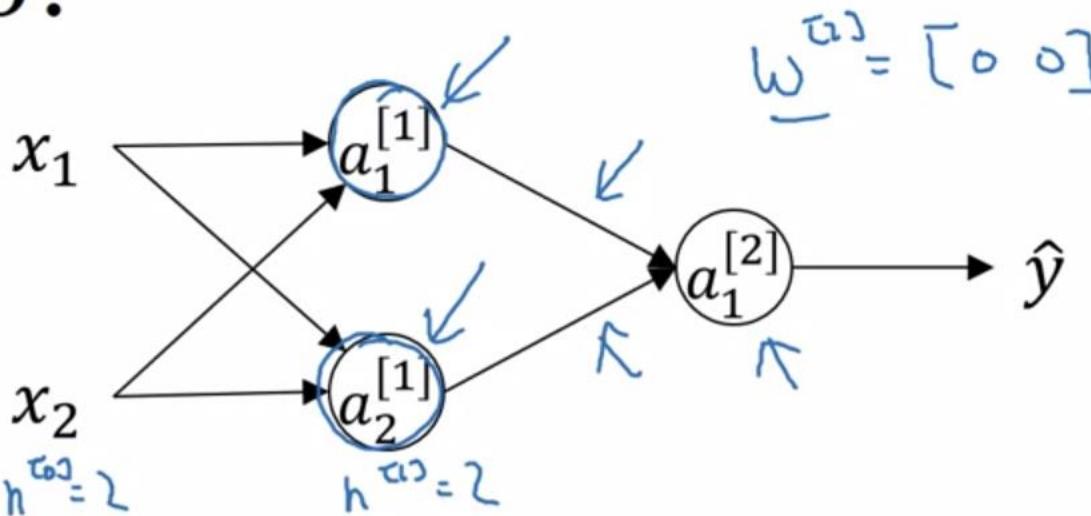
$$J = 3v$$



## Logistic regression derivatives



# What happens if you initialize weights to zero?



$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

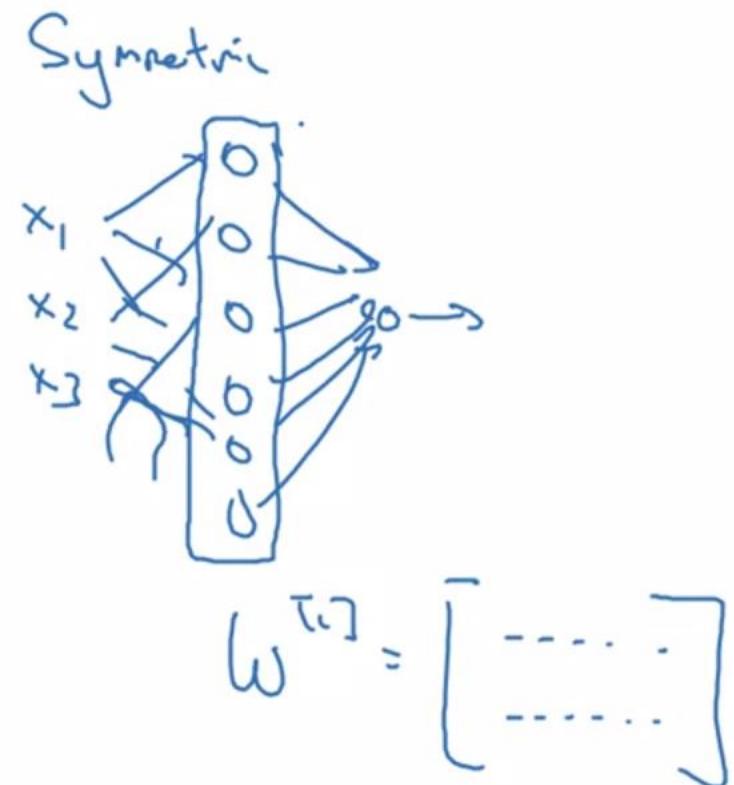
$$b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]}$$

$$\delta z_1^{[1]} = \delta z_2^{[1]}$$

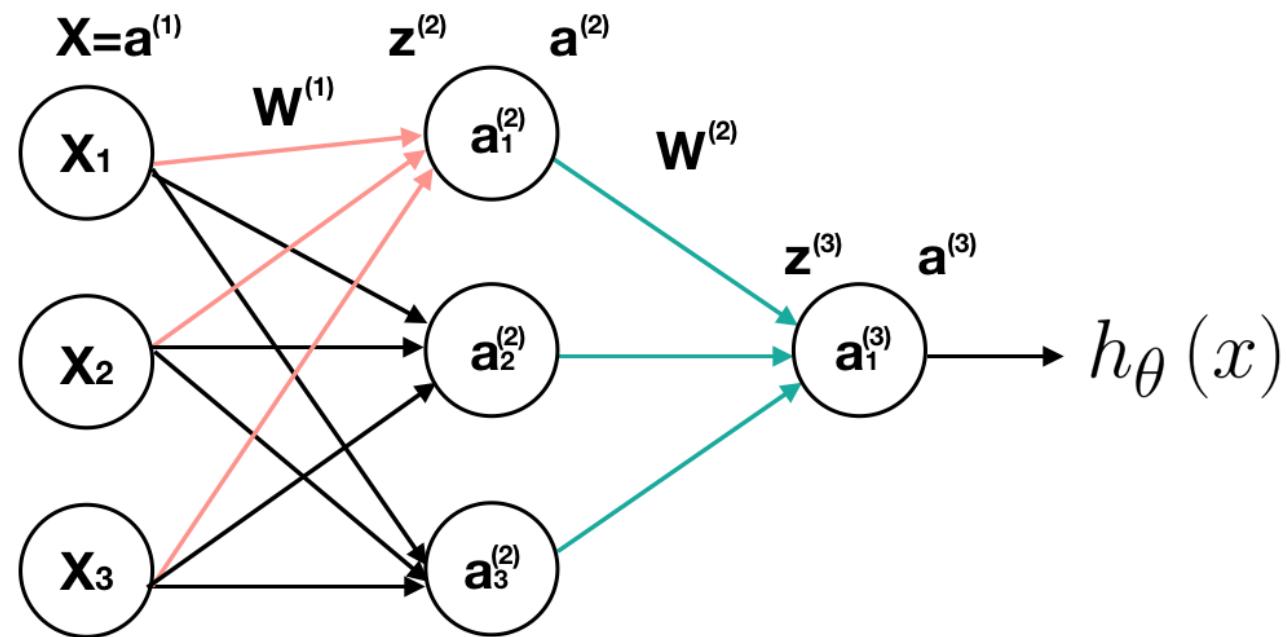
$$\delta w = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

$$w^{[1]} = w^{[1]} - \lambda \delta w$$



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

$$\Theta^{(2)} = W^{(2)} = \begin{bmatrix} \Theta_{11}^{(2)} \\ \Theta_{12}^{(2)} \\ \Theta_{13}^{(2)} \end{bmatrix} = \begin{bmatrix} w_{11}^{(2)} \\ w_{12}^{(2)} \\ w_{13}^{(2)} \end{bmatrix}$$



$$\frac{\partial J(\Theta)}{\Theta^{(2)}} = \delta^{(3)} * a^{(2)}$$

$$\frac{\partial J(\Theta)}{\Theta^{(1)}} = \delta^{(2)} * a^{(1)}$$

$$\delta^{(3)} = (a^{(3)} - y) * g'(z^{(3)})$$

$$\delta^{(2)} = \delta^{(3)} (W^{(2)})^T g'(z^{(2)})$$

**Hint**

$$\delta^{(3)} = (a^{(3)} - y) * g(z^{(3)})$$

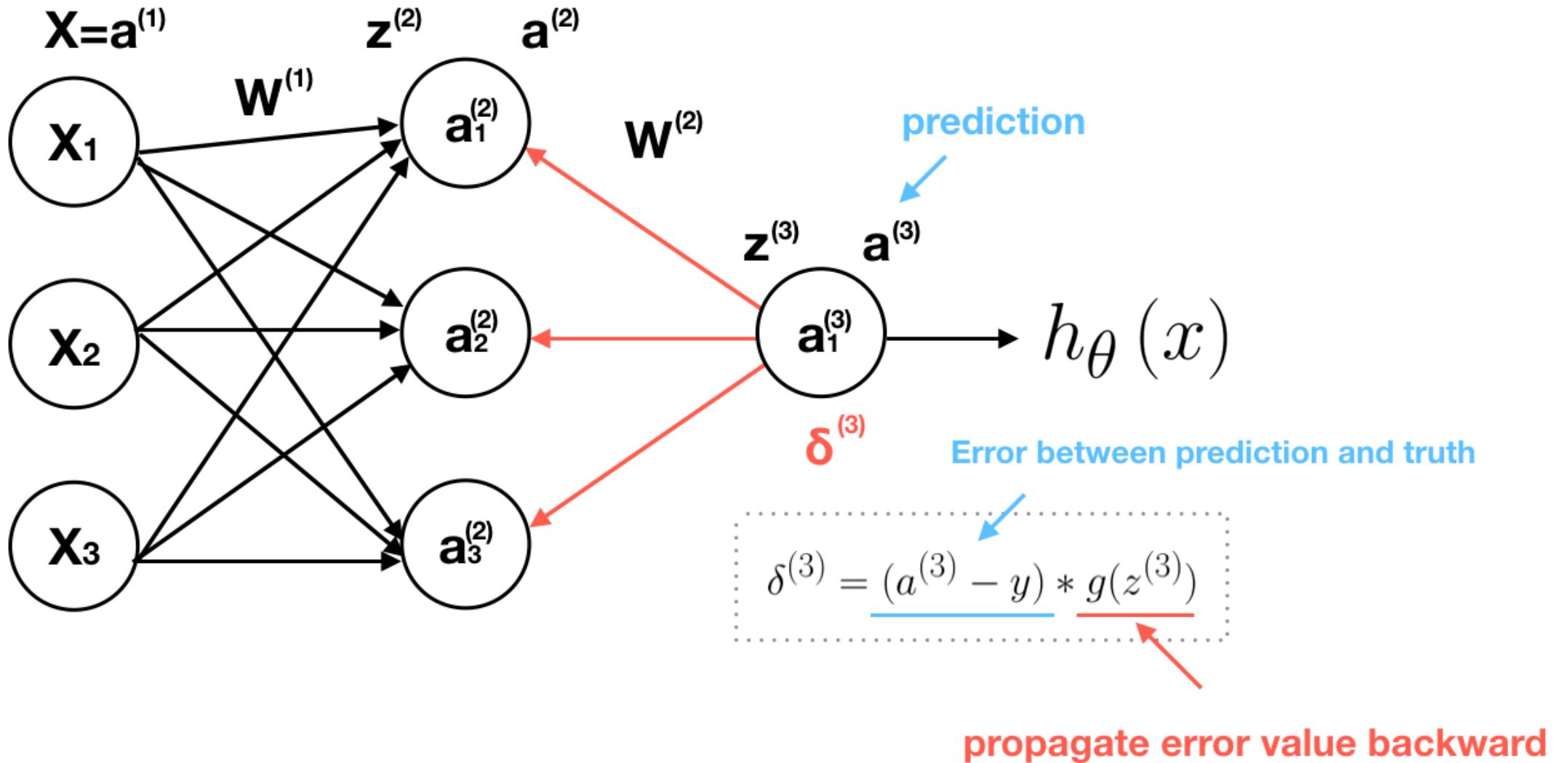
$$\begin{bmatrix} \delta_1^{(3)} \\ \delta_2^{(3)} \\ \delta_3^{(3)} \end{bmatrix} = \begin{bmatrix} (a_1^{(3)} - y_1) \\ (a_2^{(3)} - y_2) \\ (a_3^{(3)} - y_3) \end{bmatrix} * \begin{bmatrix} g'(z_1^{(3)}) \\ g'(z_2^{(3)}) \\ g'(z_3^{(3)}) \end{bmatrix} = \begin{bmatrix} (a_1^{(3)} - y_1) * g'(z_1^{(3)}) \\ (a_2^{(3)} - y_2) * g'(z_2^{(3)}) \\ (a_3^{(3)} - y_3) * g'(z_3^{(3)}) \end{bmatrix}$$

**Hint**

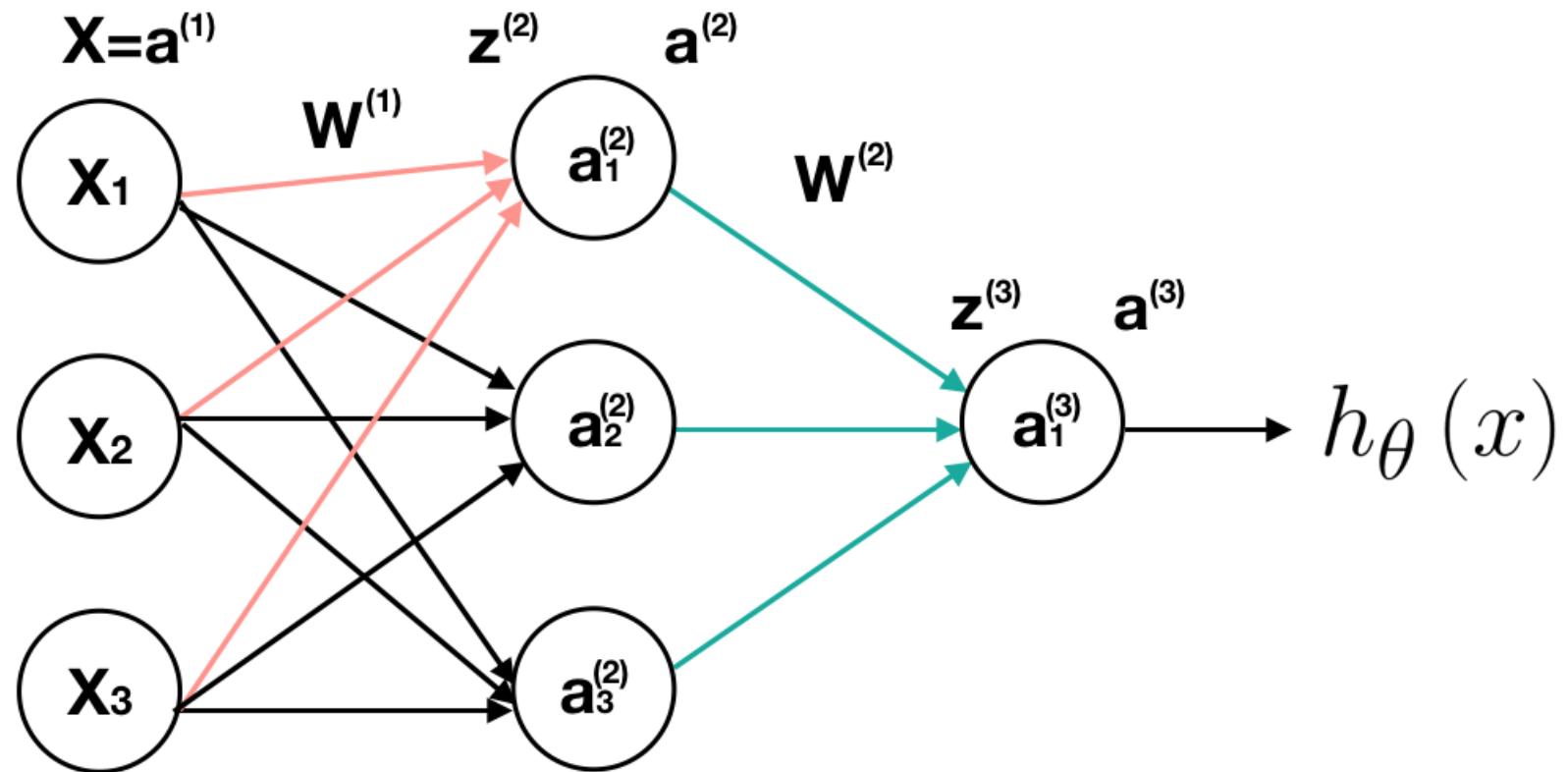
$$\delta^{(2)} = \delta^{(3)} (W^{(2)})^T g'(z^{(2)})$$

$$= \left( \begin{bmatrix} \delta_1^{(3)} \\ \delta_2^{(3)} \\ \delta_3^{(3)} \end{bmatrix} [w_1^{(2)} \ w_2^{(2)} * w_3^{(2)}] \right) * \begin{bmatrix} g'(z_1^{(2)}) \\ g'(z_2^{(2)}) \\ g'(z_3^{(2)}) \end{bmatrix}$$

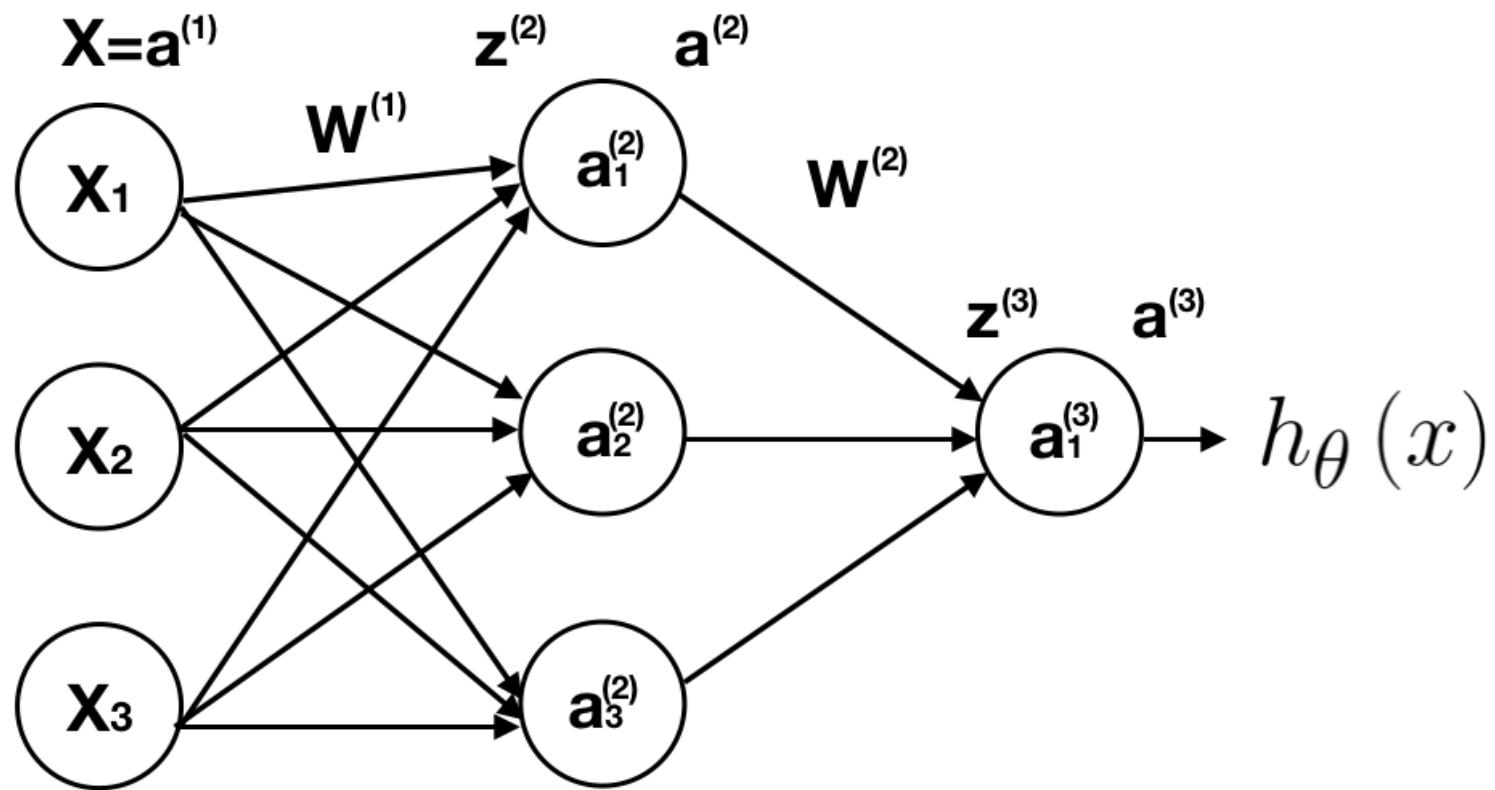
$$= \begin{bmatrix} \delta_1^{(3)} * w_1^{(2)} & \delta_1^{(3)} * w_2^{(2)} & \delta_1^{(3)} * w_3^{(2)} \\ \delta_2^{(3)} * w_2^{(2)} & \delta_2^{(3)} * w_2^{(2)} & \delta_2^{(3)} * w_3^{(2)} \\ \delta_3^{(3)} * w_3^{(2)} & \delta_3^{(3)} * w_2^{(2)} & \delta_3^{(3)} * w_3^{(2)} \end{bmatrix} * \begin{bmatrix} g'(z_1^{(2)}) \\ g'(z_2^{(2)}) \\ g'(z_3^{(2)}) \end{bmatrix}$$



$$\Theta^{(1)} = W^{(1)} = \boxed{\begin{bmatrix} \Theta_{11}^{(1)} & \Theta_{21}^{(1)} & \Theta_{31}^{(1)} \\ \Theta_{12}^{(1)} & \Theta_{22}^{(1)} & \Theta_{32}^{(1)} \\ \Theta_{13}^{(1)} & \Theta_{23}^{(1)} & \Theta_{33}^{(1)} \end{bmatrix}} \quad \Theta^{(2)} = W^{(2)} = \boxed{\begin{bmatrix} \Theta_{11}^{(2)} \\ \Theta_{12}^{(2)} \\ \Theta_{13}^{(2)} \end{bmatrix}}$$



$$J(\Theta) = -[y \log(h_\Theta(x)) + (1-y) \log(1-h_\Theta(x))]$$



### Forward Propagation

$$(\Theta^{(1)})^T X = z^{(2)}$$

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

$$(\Theta^{(2)})^T a^{(2)} = z^{(3)}$$

$$a^{(3)} = g(z^{(3)})$$

**Hint**

$$a^{(3)} = g(z^{(3)})$$

$$\frac{\partial J(\Theta)}{\Theta_{11}^{(2)}} = \boxed{\frac{\partial J(\Theta)}{\partial a_1^{(3)}}} * \boxed{\frac{\partial a_1^{(3)}}{\partial z_1^{(3)}}} * \boxed{\frac{\partial z_1^{(3)}}{\Theta_{11}^{(2)}}} = \underline{\delta_1^{(3)} * a_1^{(2)}}$$

\$a^{(3)} - y\$  
 \$g'(z^{(3)})\$  
 \$a\_1^{(2)}\$

**Hint**

$$\Theta^{(2)} a^{(2)} = z^{(3)}$$

**Hint**

$$\begin{aligned} \frac{\partial J(\Theta)}{a_1^{(3)}} &= -[y * \frac{1}{g(z)} * g'(z) + (1 - y) * \frac{1}{1 - g(z)} * (-g'(z))] \\ &= -[y * (1 - g(z)) + (y - 1) * g(z)] = g(z) - y \end{aligned}$$

**Hint-The Property of Sigmoid Function**

$$g'(z) = g(z) * (1 - g(z))$$

$$\frac{\partial J(\Theta)}{\Theta_{11}^{(1)}} = \boxed{\frac{\partial J(\Theta)}{\partial a_1^{(3)}} * \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}}} * \boxed{\frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} * \boxed{\frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} * \boxed{\frac{\partial z_1^{(2)}}{\Theta_{11}^{(1)}}}}} = \underline{\delta_1^{(2)} * a_1^{(1)}}$$

**Hint**

$$(\Theta^{(1)})^T X = z^{(2)}$$

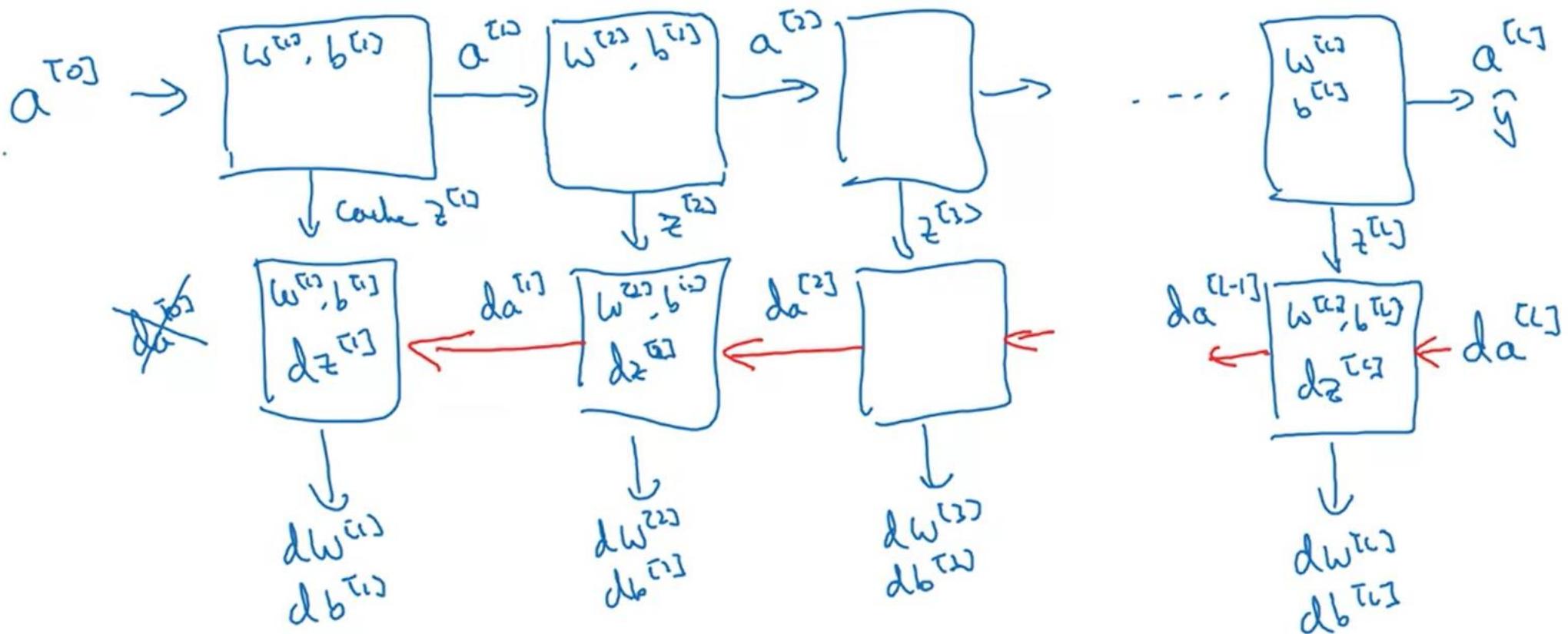
**Hint**

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

**Hint**

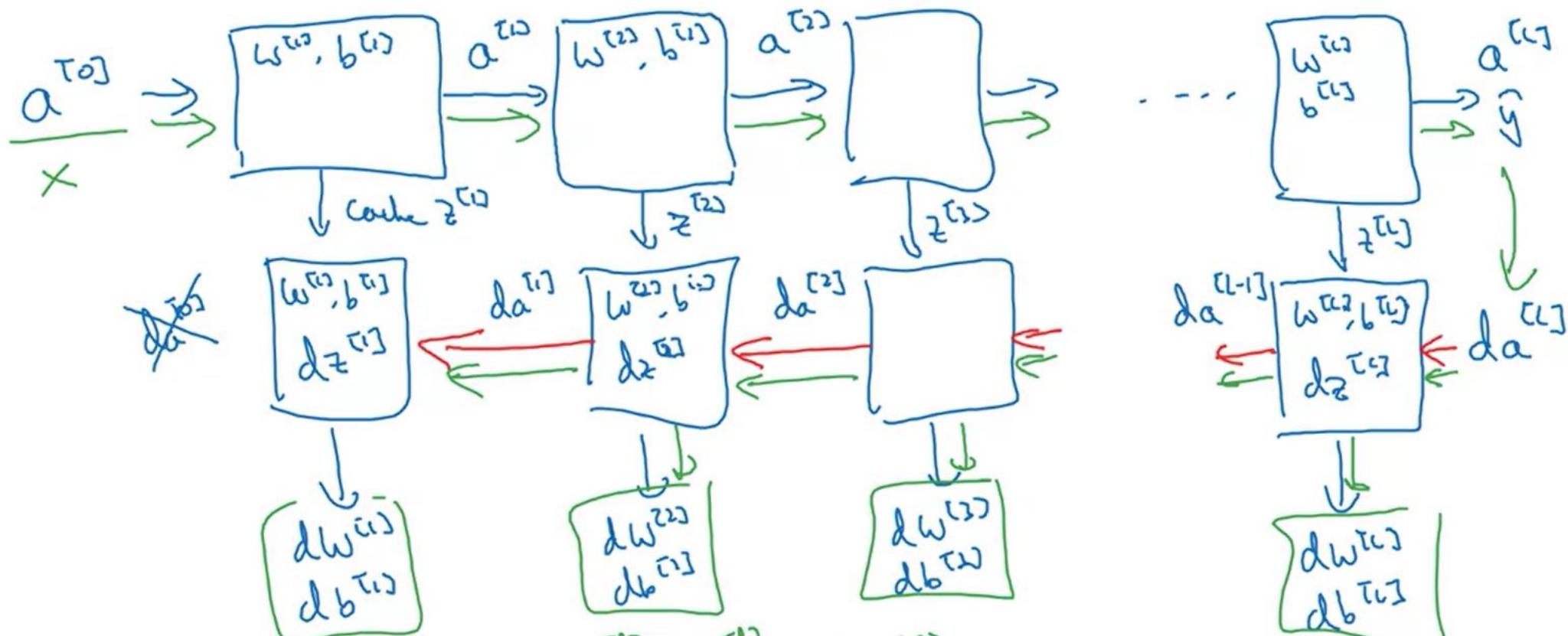
$$\Theta^{(2)} a^{(2)} = z^{(3)}$$

# Forward and backward functions



Andrew Ng

# Forward and backward functions



$$w^{(l)} := w^{(l)} - \alpha dw^{(l)}$$

$$b^{(l)} := b^{(l)} - \alpha db^{(l)}$$

Andrew Ng

# Forward propagation for layer $l$

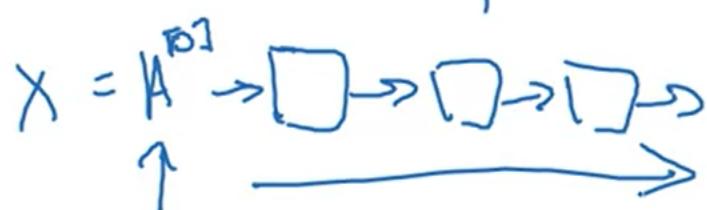
→ Input  $a^{[l-1]} \leftarrow$

→ Output  $a^{[l]}$ , cache ( $\underline{z^{[l]}}$ )

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

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

$$\begin{matrix} a^{[0]} \\ A^{[0]} \end{matrix}$$



Vertwijf:

$$z^{[l]} = w^{[l]} \cdot A^{[l-1]} + b^{[l]}$$

$$A^{[l]} = g^{[l]}(z^{[l]})$$

# Computing gradients

## Logistic regression

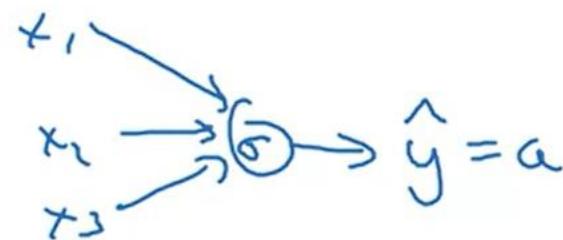


Diagram illustrating the computation of gradients for a logistic regression model:

The forward pass consists of:

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

where  $\sigma(z) = g(z) = \frac{1}{1+e^{-z}}$ .

The backward pass (gradients) are computed as follows:

$$\begin{aligned} dz &= a - y \\ dz &= da \cdot g'(z) \\ \frac{\partial \mathcal{L}}{\partial z} &= \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{da}{dz} \end{aligned}$$

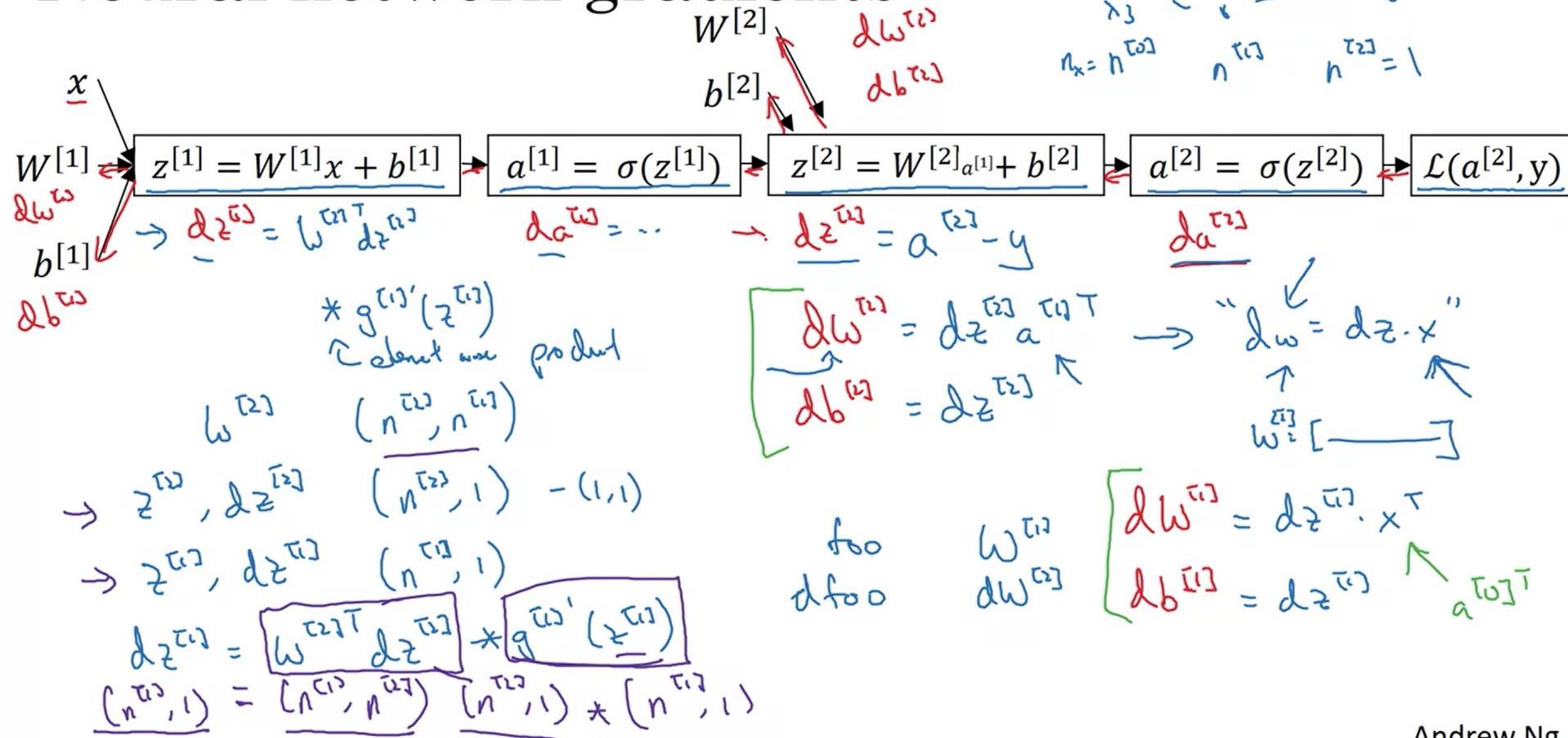
Specifically, the gradient of the loss function with respect to  $a$  is:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial a} &= \frac{d}{da} \mathcal{L}(a, y) = -y \log a - (1-y) \log(1-a) \\ &= -\frac{y}{a} + \frac{1-y}{1-a} \end{aligned}$$

And the gradient of the sigmoid function with respect to  $z$  is:

$$\frac{d}{dz} g(z) = g'(z)$$

# Neural network gradients



# Summary of gradient descent

$$\underline{dz}^{[2]} = \underline{a}^{[2]} - \underline{y}$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

(n<sup>T<sub>1</sub></sup>, 1)

$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

$$\underline{dZ}^{[2]} = \underline{A}^{[2]} - \underline{Y}$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis = 1, keepdims = True)$$

$$dZ^{[1]} = \underbrace{W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]})}_{\substack{\text{elementwise product} \\ (n^{T_1}, m) \quad (n^{T_2}, m) \quad (n^{T_1}, m)}}$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True)$$

$$J(\cdot) = \frac{1}{m} \sum_{i=1}^n \mathcal{L}(\hat{y}_i, y_i)$$