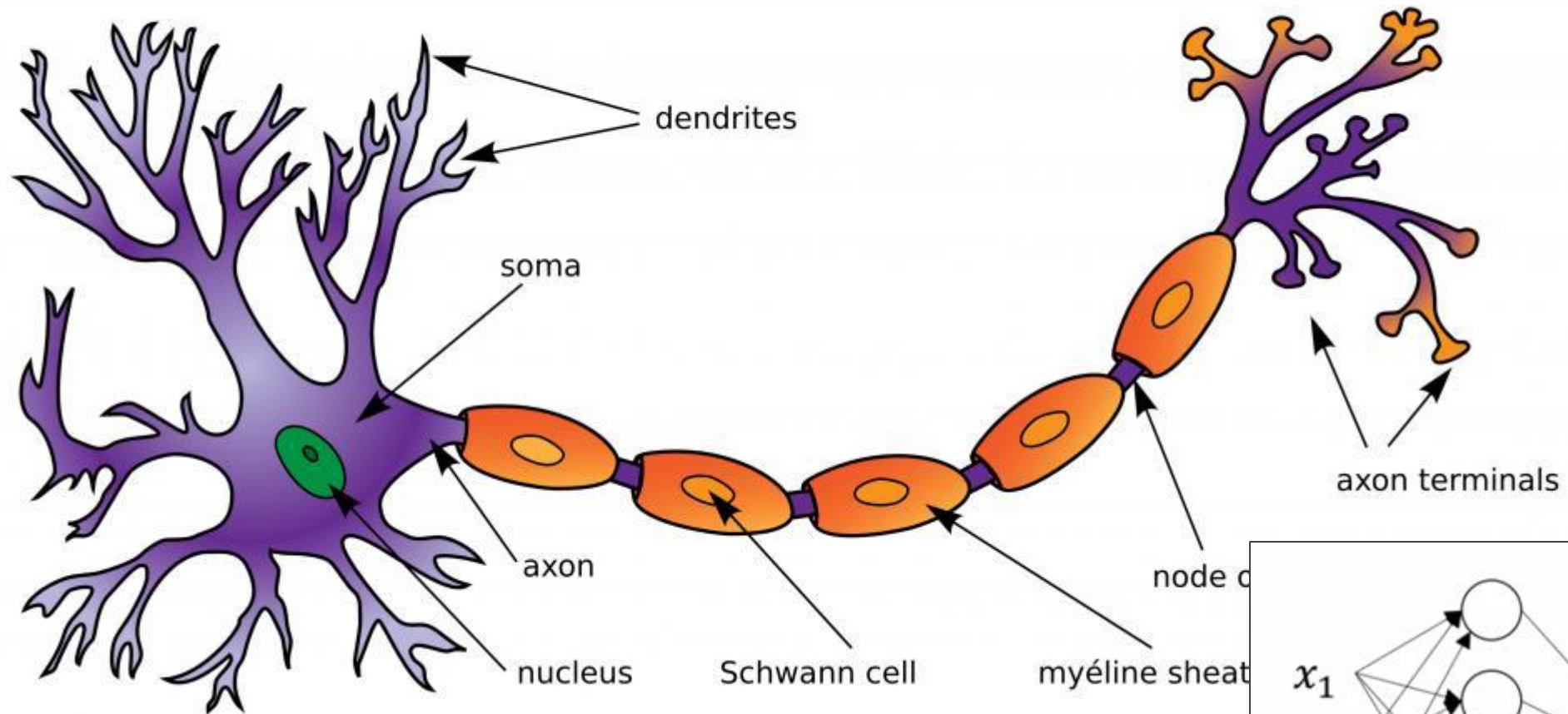# Shallow Neural Networks

= 1 hidden layer
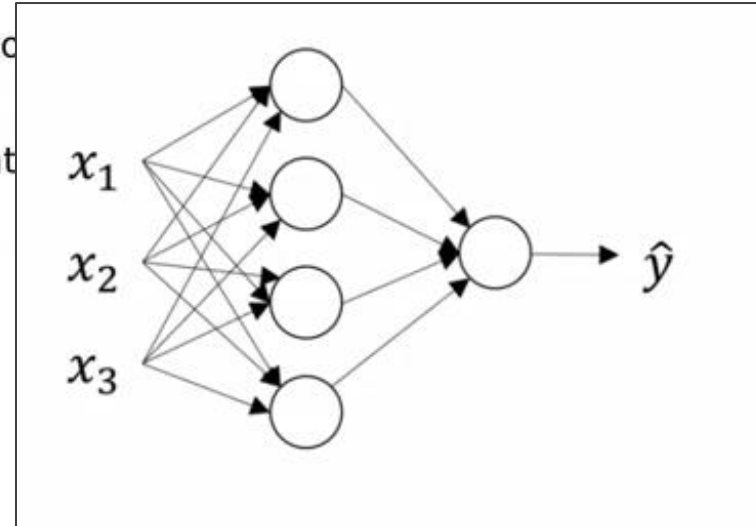
2021년 3월 20일 토요일 - 발표자 : 최윤정

20 Mar 2021 Sat- Speaker : Yoon Choi

# NEURON
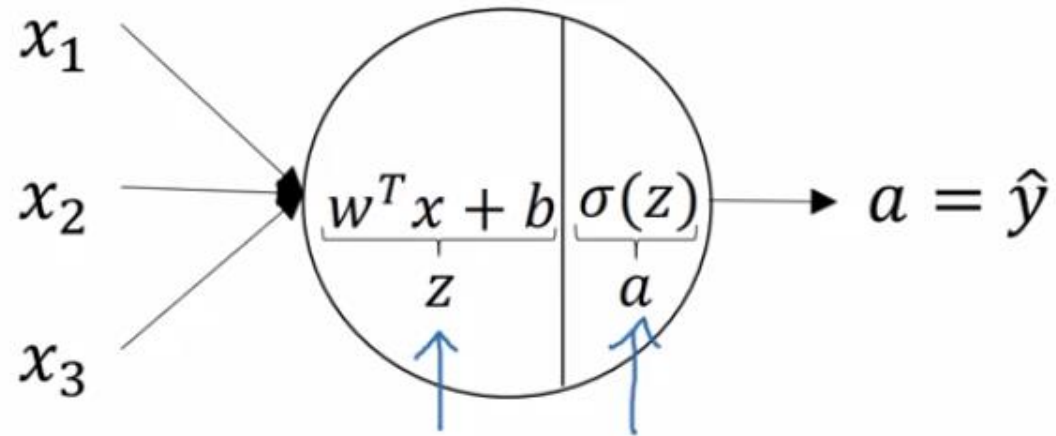


dendrites

soma

axon terminals
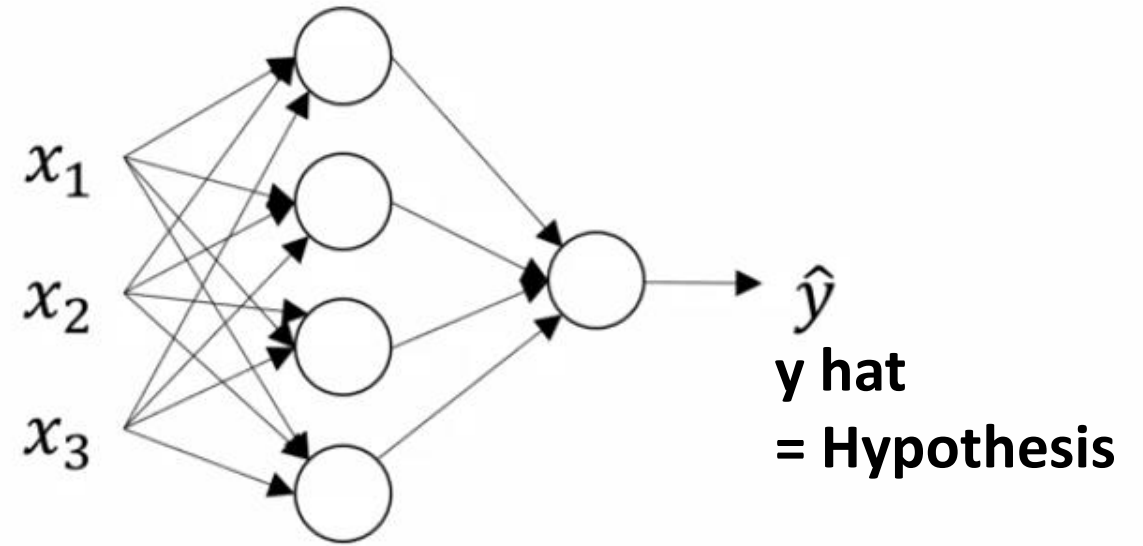
axon

node o

nucleus      Schwann cell      myéline sheat

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

$$a = \sigma(z)$$

$x_1$

$x_2$

$x_3$

$\hat{y}$

# Neural Network Representation



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

$$a = \sigma(z)$$

y hat
= Hypothesis

# Neural Network Representation

$$z_1^{[1]} = w_1^{[1]T}x + b_1^{[1]}$$
$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$$a_i^{[\ell]} \leftarrow \text{layer}$$
$$a_i \leftarrow \text{node in layer}$$

$x_1$

$x_2$

$x_3$

$$w^T x + b \quad \sigma(z) \longrightarrow a = \hat{y}$$
$$\underbrace{\qquad}_{z} \quad \underbrace{\qquad}_{a}$$

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

$$a = \sigma(z)$$

$x_1$

$x_2$

$x_3$

$\hat{y}$

$$z_2^{[1]} = w_2^{[1]T}x + b_2^{[1]}$$
$$a_2^{[1]} = \sigma(z_2^{[1]})$$

$x_1$

$x_2$

$x_3$

NEURON

dendrites

soma

axon terminals

axon

node of Ranvier

nucleus   Schwann cell   myéline sheath

**cost function of Classification**
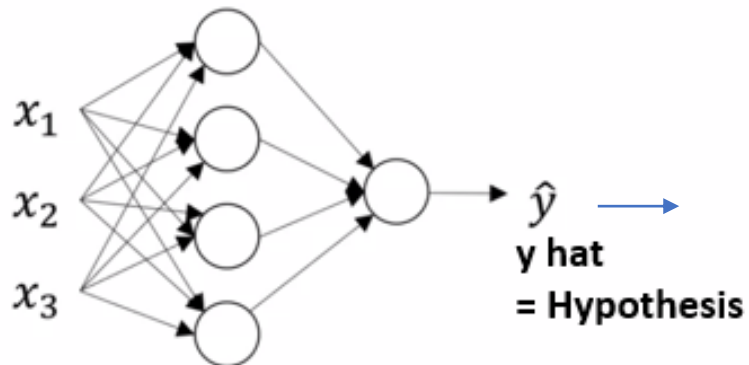
$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{m} [\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))]$$

**cost function of Neural Network**

$$\mathcal{L}^{(i)}(\hat{y}^{(i)}, y^{(i)}) = -(y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$$



$x_1$

$x_2$

$x_3$

$\hat{y}$ $\longrightarrow$

y hat

= Hypothesis

$$\mathcal{L}^{(i)}(\hat{y}^{(i)}, y^{(i)}) = -(y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$$

# Gradient descent for neural networks

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$
$(n^{[1]}, n^{[0]})$ $(n^{[1]}, 1)$ $(n^{[2]}, n^{[1]})$ $(n^{[2]}, 1)$

$n_x = n^{[0]}, n^{[1]}, n^{[2]} = 1$

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

$\uparrow a^{[2]}$

Gradient descent:

**Gradient of the Cost function을 구하기위해**
**Back propagation 사용**

$\rightarrow$ Repeat $\{$

Compute predicts $(\hat{y}^{(i)}, i=1,\dots,m)$

$dW^{[i]} = \frac{\partial J}{\partial W^{[i]}}, \quad db^{[i]} = \frac{\partial J}{\partial b^{[i]}}, \dots$

$W^{[i]} := W^{[i]} - \alpha \, dW^{[i]}$

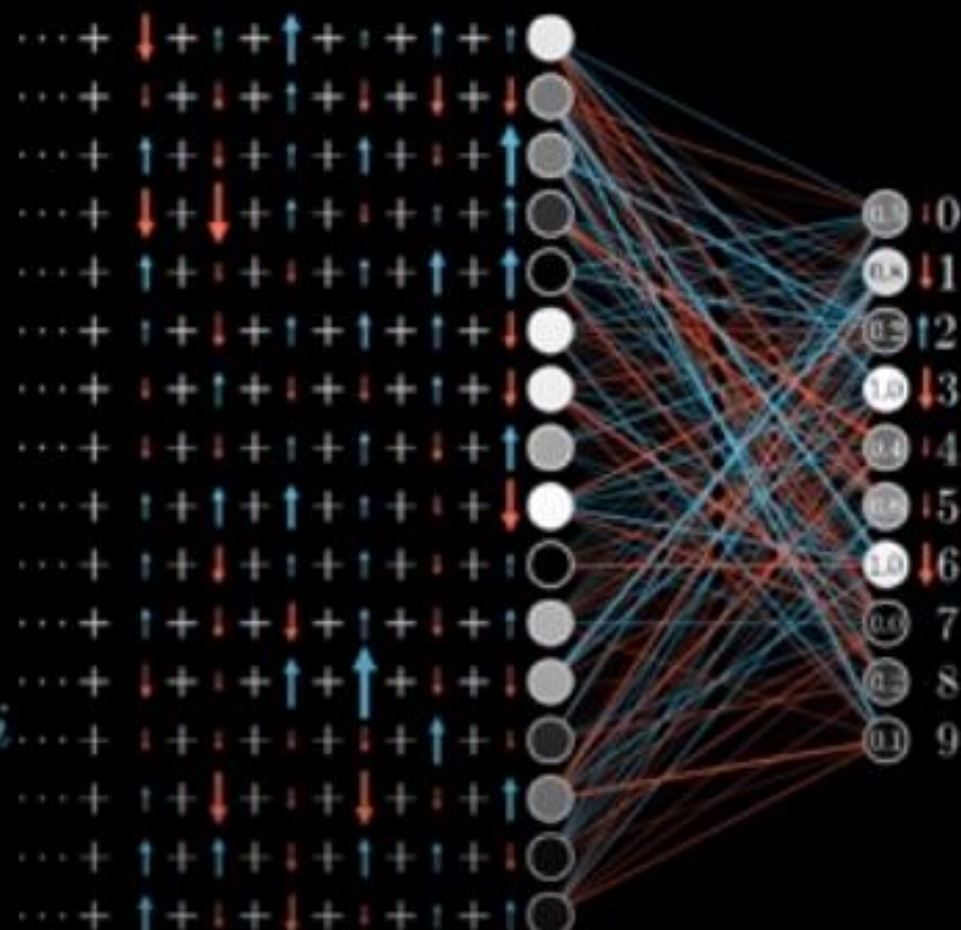$b^{[i]} := b^{[i]} - \alpha \, db^{[i]}$

Propagate backwards

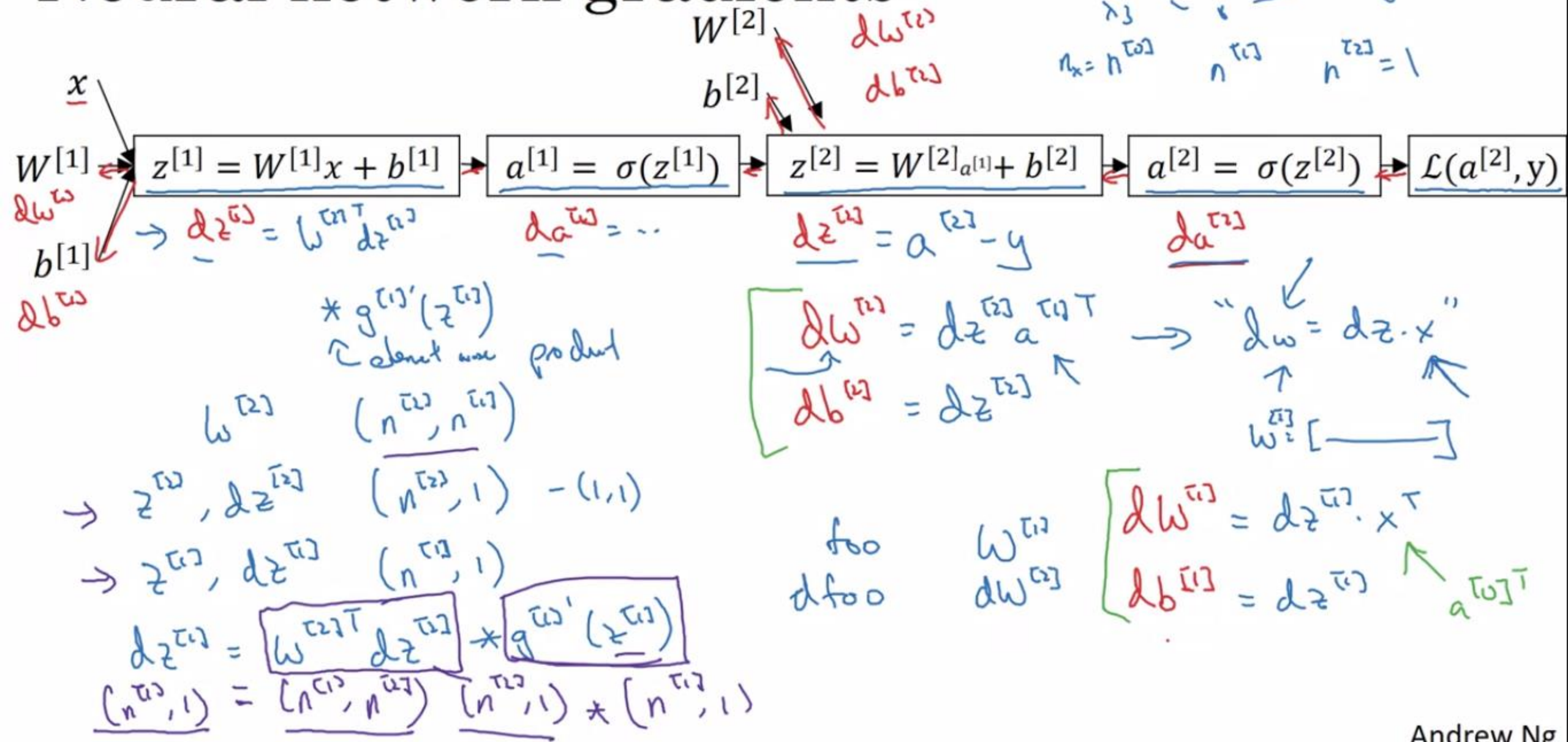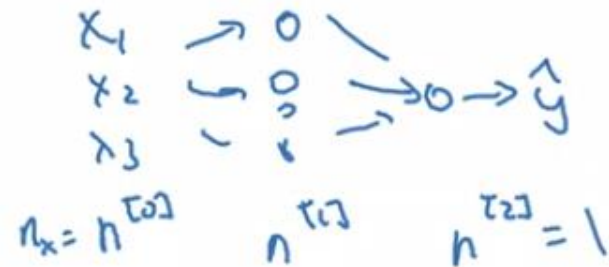$$\mathcal{L}^{(i)}(\hat{y}^{(i)}, y^{(i)})$$

Increase $b$

Increase $w_i$
in proportion to $a_i$

Change $a_i$
in proportion to $w_i$

# Neural network gradients



$W^{[2]}$   $dW^{[2]}$   $db^{[2]}$

$b^{[2]}$

$x$

$W^{[1]}$   $dW^{[1]}$

$b^{[1]}$   $db^{[1]}$

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

$\rightarrow dz^{[1]} = W^{[2]T} dz^{[2]}$

$da^{[1]} = \cdots$

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

$da^{[2]}$

$x_1 \rightarrow 0$
$x_2 \rightarrow 0 \rightarrow 0 \rightarrow \hat{y}$
$x_3$

$n_x = n^{[0]}$   $n^{[1]}$   $n^{[2]} = 1$

$* \; g^{[1]'}(z^{[1]})$
element wise product

$W^{[2]}$   $(n^{[2]}, n^{[1]})$

$\rightarrow z^{[2]}, dz^{[2]} \quad (n^{[2]}, 1) - (1,1)$

$\rightarrow z^{[1]}, dz^{[1]} \quad (n^{[1]}, 1)$

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

$(n^{[1]}, 1) = (n^{[1]}, n^{[2]}) \; (n^{[2]}, 1) * (n^{[1]}, 1)$

$\begin{cases} dW^{[2]} = dz^{[2]} a^{[1]T} \\ db^{[2]} = dz^{[2]} \end{cases} \rightarrow \text{``} dw = dz \cdot x \text{''}$

$W^{[2]}: [\underline{\quad\quad}]$

$foo \quad W^{[1]}$
$dfoo \quad dW^{[2]}$

$\begin{cases} dW^{[1]} = dz^{[1]} \cdot x^T \\ db^{[1]} = dz^{[1]} \end{cases} \quad a^{[0]T}$

Andrew Ng

# Summary of gradient descent

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

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

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

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

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

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

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

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

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

elementwise product

$$dZ^{[1]} = \underbrace{W^{[2]^T} dZ^{[2]}}_{(n^{[1]}, m)} * \underbrace{g^{[1]'}(Z^{[1]})}_{(n^{[1]}, m)}$$
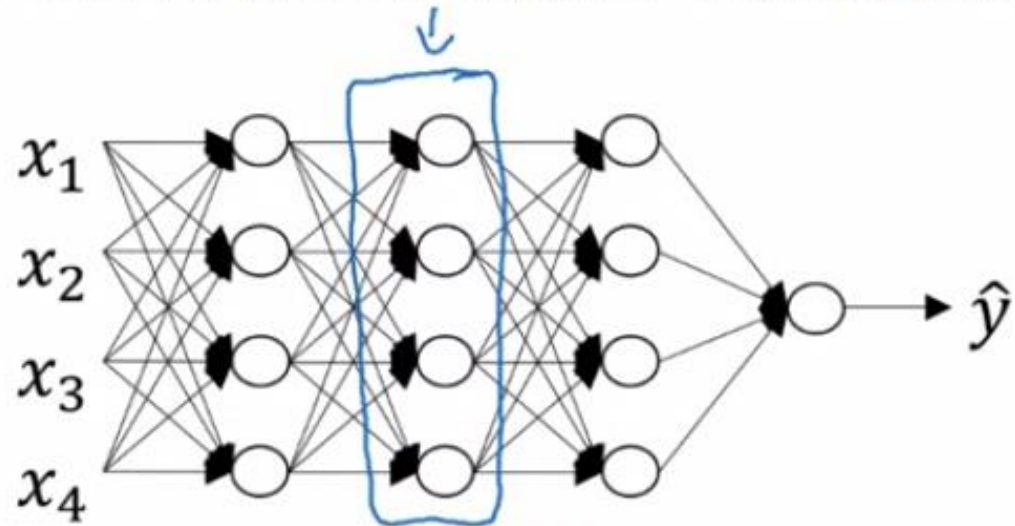$$(n^{[1]}, m)$$

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

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

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

# Deep Neural Networks

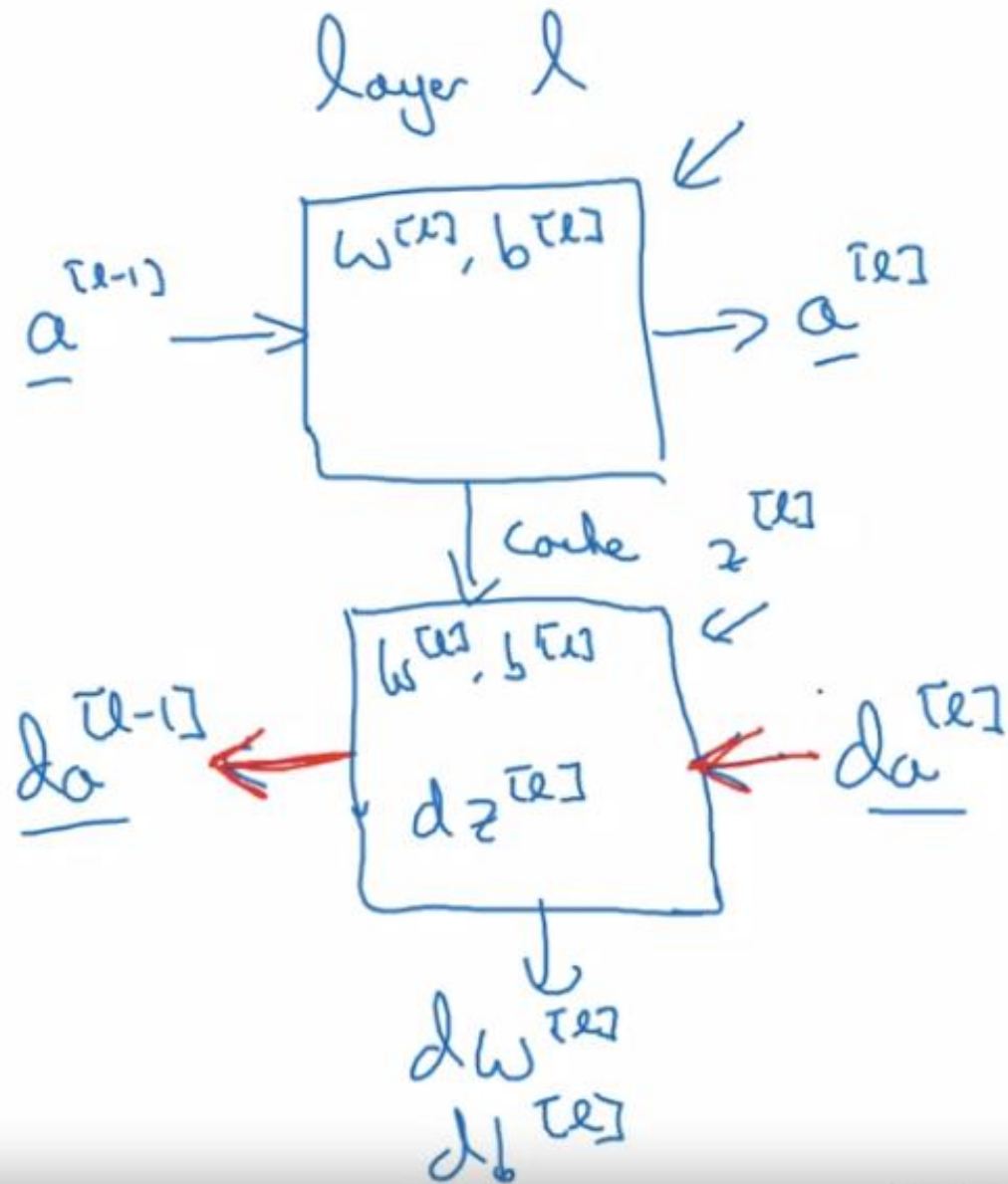= more than 1 hidden layer

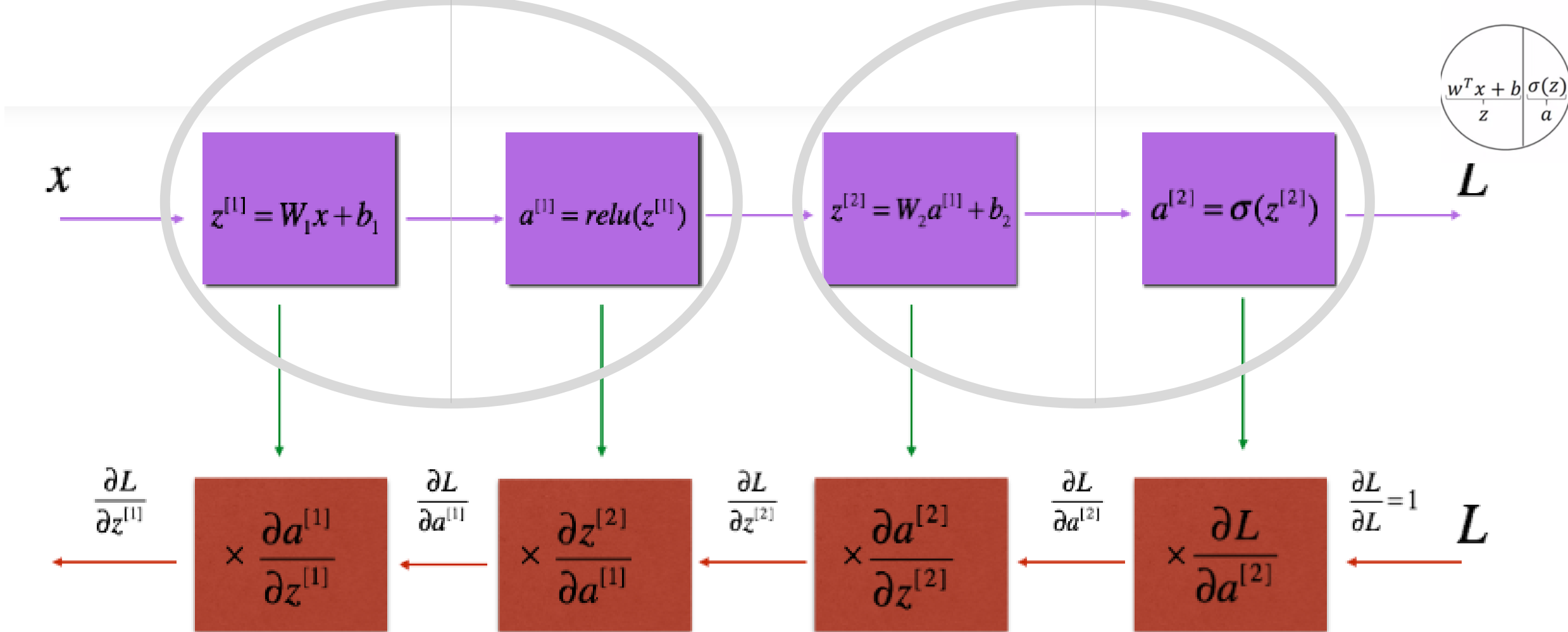# Forward and backward functions



Layer $l$: $W^{[l]}$, $b^{[l]}$

→ Forward: Input $a^{[l-1]}$, output $a^{[l]}$

$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$    cache $z^{[l]}$

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

→ Backward: Input $da^{[l]}$, output $da^{[l-1]}$

cache $(z^{[l]})$    $dW^{[l]}$

$db^{[l]}$

Andrew Ng

Now, similar to forward propagation, you are going to build the backward propagation in three steps:

- LINEAR backward
- LINEAR -> ACTIVATION backward where ACTIVATION computes the derivative of either the ReLU or sigmoid activation
- [LINEAR -> RELU] × (L-1) -> LINEAR -> SIGMOID backward (whole model)

# Backward propagation for layer $l$

$$dW^{[l]} = \frac{\partial J}{\partial W^{[l]}} = \frac{1}{m} dZ^{[l]} A^{[l-1]T}$$

$$db^{[l]} = \frac{\partial J}{\partial b^{[l]}} = \frac{1}{m} \sum_{i=1}^{m} dZ^{[l](i)}$$

$$dA^{[l-1]} = \frac{\partial \mathcal{L}}{\partial A^{[l-1]}} = W^{[l]T} dZ^{[l]}$$

→ Input $da^{[l]}$

→ Output $da^{[l-1]}, dW^{[l]}, db^{[l]}$

$$dz^{[l]} = da^{[l]} * g^{[l]'}(z^{[l]})$$

$$dW^{[l]} = dz^{[l]} \cdot a^{[l-1]}$$

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

$$da^{[l-1]} = W^{[l]T} \cdot dz^{[l]}$$

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

$$dz^{[l]} = dA^{[l]} * g^{[l]'}(z^{[l]})$$

$$dW^{[l]} = \frac{1}{m} dz^{[l]} \cdot A^{[l-1]T}$$

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

$$dA^{[l-1]} = W^{[l]T} \cdot dz^{[l]}$$

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

# What are hyperparameters?

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, W^{[3]}, b^{[3]}$ ...

Hyperparameters: Learning rate $\alpha$

#iterations

#hidden layers $L$
#hidden units $n^{[1]}, n^{[2]}, \dots$
choice of activation function

Later: Momentum, mini-batch size, regularizations, ...

**hyperparameter : parameter의 값이나 상태를 결정짓는 요소들**

Andrew Ng

- [x] weight matrices $W^{[l]}$

  > ❗ **This should not be selected**

- [x] bias vectors $b^{[l]}$

  > ❗ **This should not be selected**

- [x] number of layers $L$ in the neural network

  > ✔ **Correct**

- [x] size of the hidden layers $n^{[l]}$

  > ✔ **Correct**

- [x] learning rate $\alpha$

  > ✔ **Correct**

- [x] number of iterations

  > ✔ **Correct**

- [x] activation values $a^{[l]}$

  > ❗ **This should not be selected**