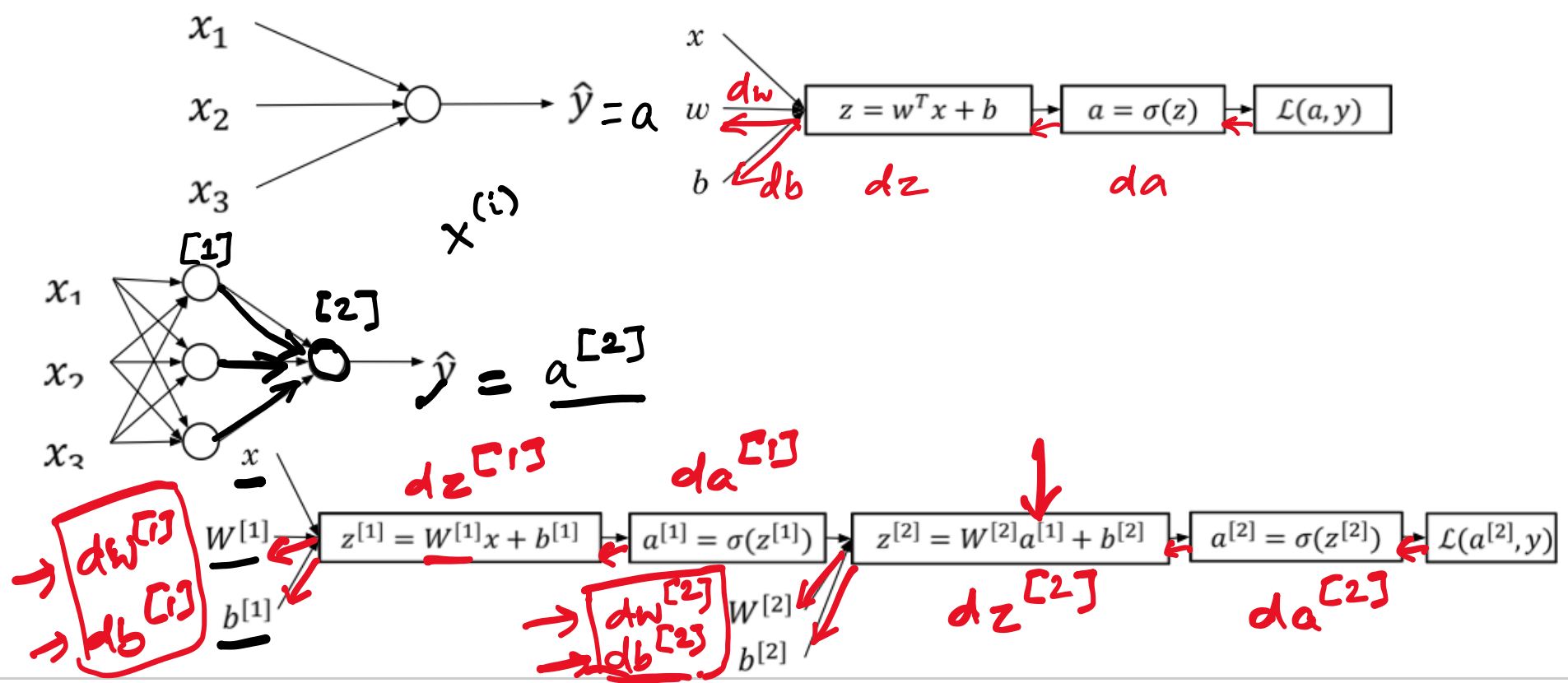
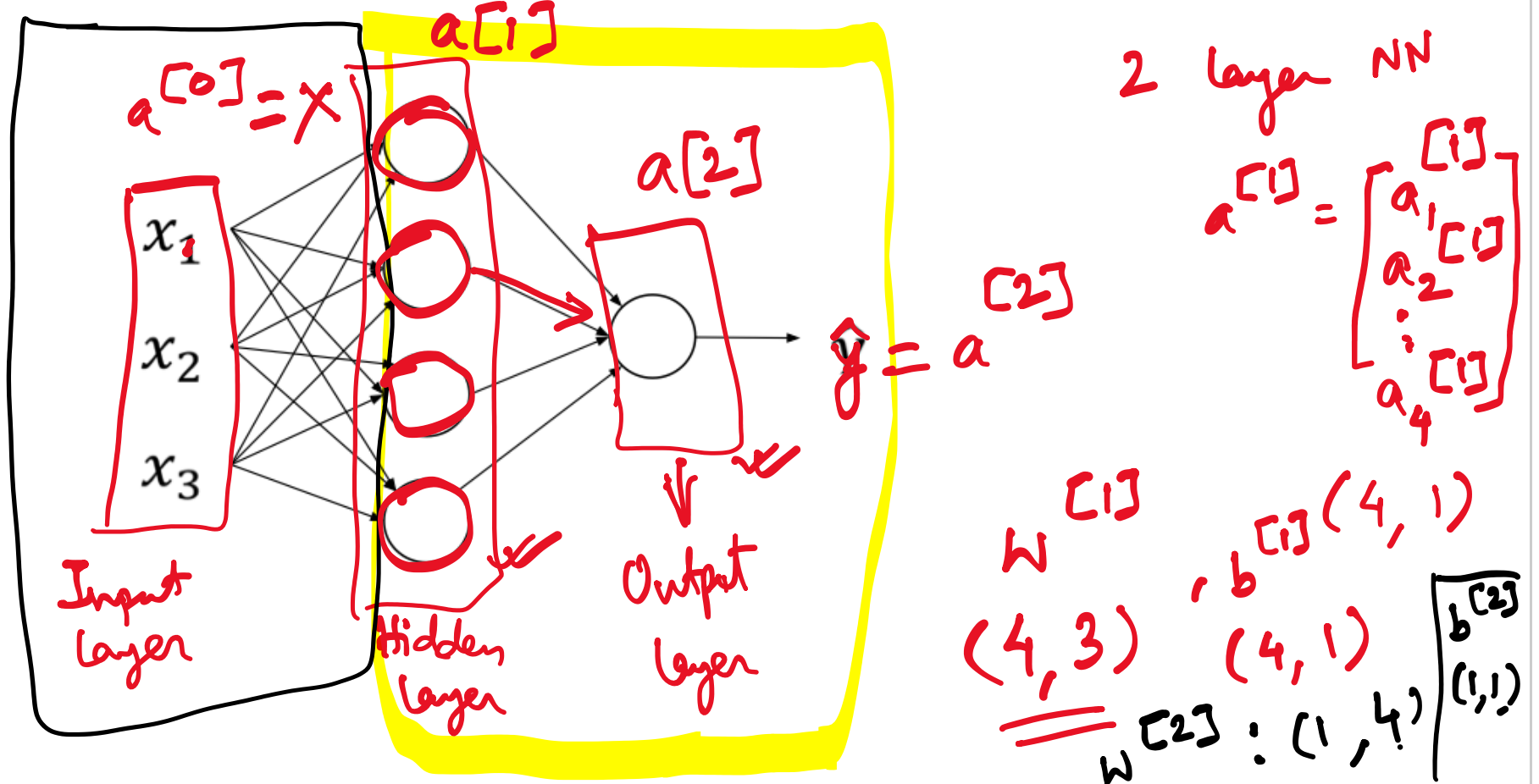


Neural Networks & Deep Learning - Part II

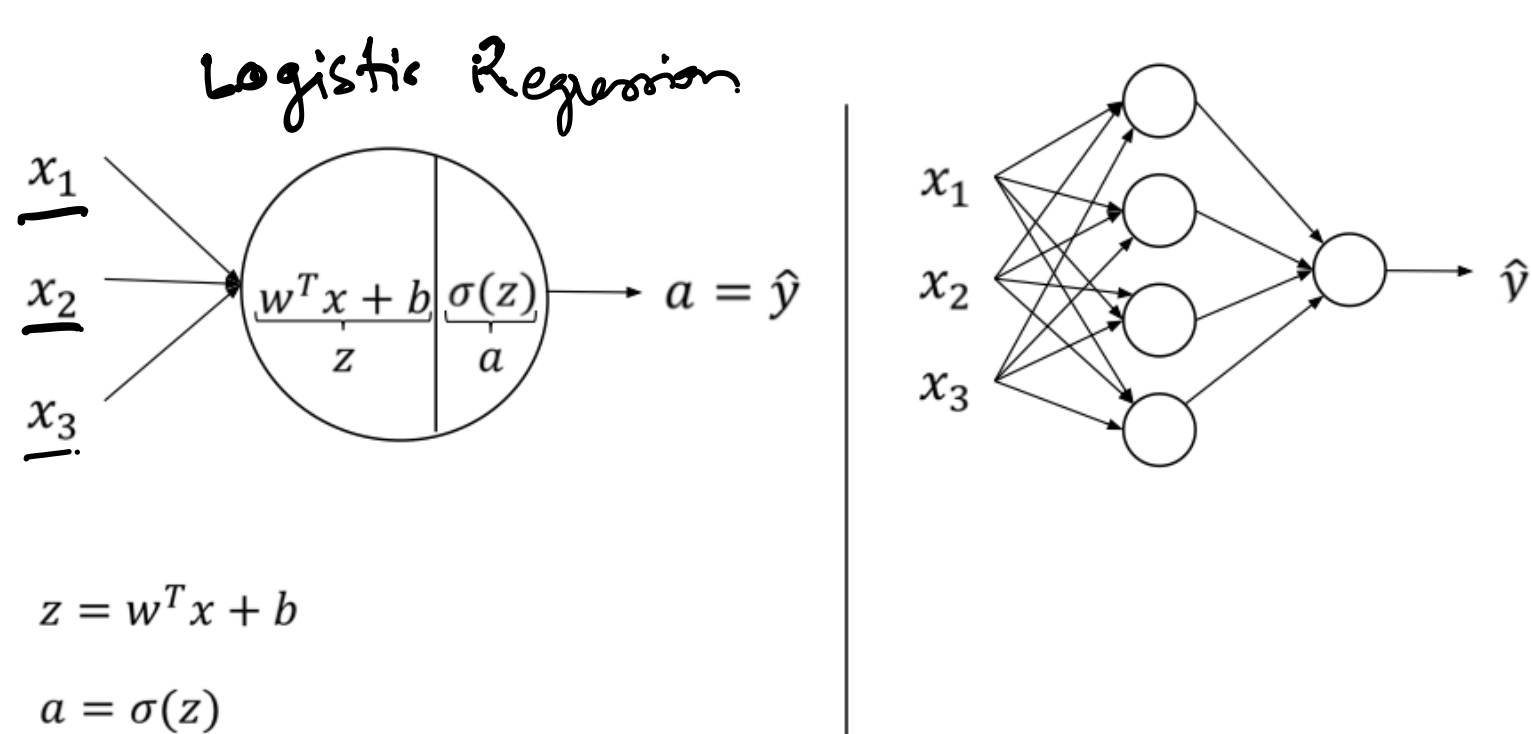
Neural Networks Overview



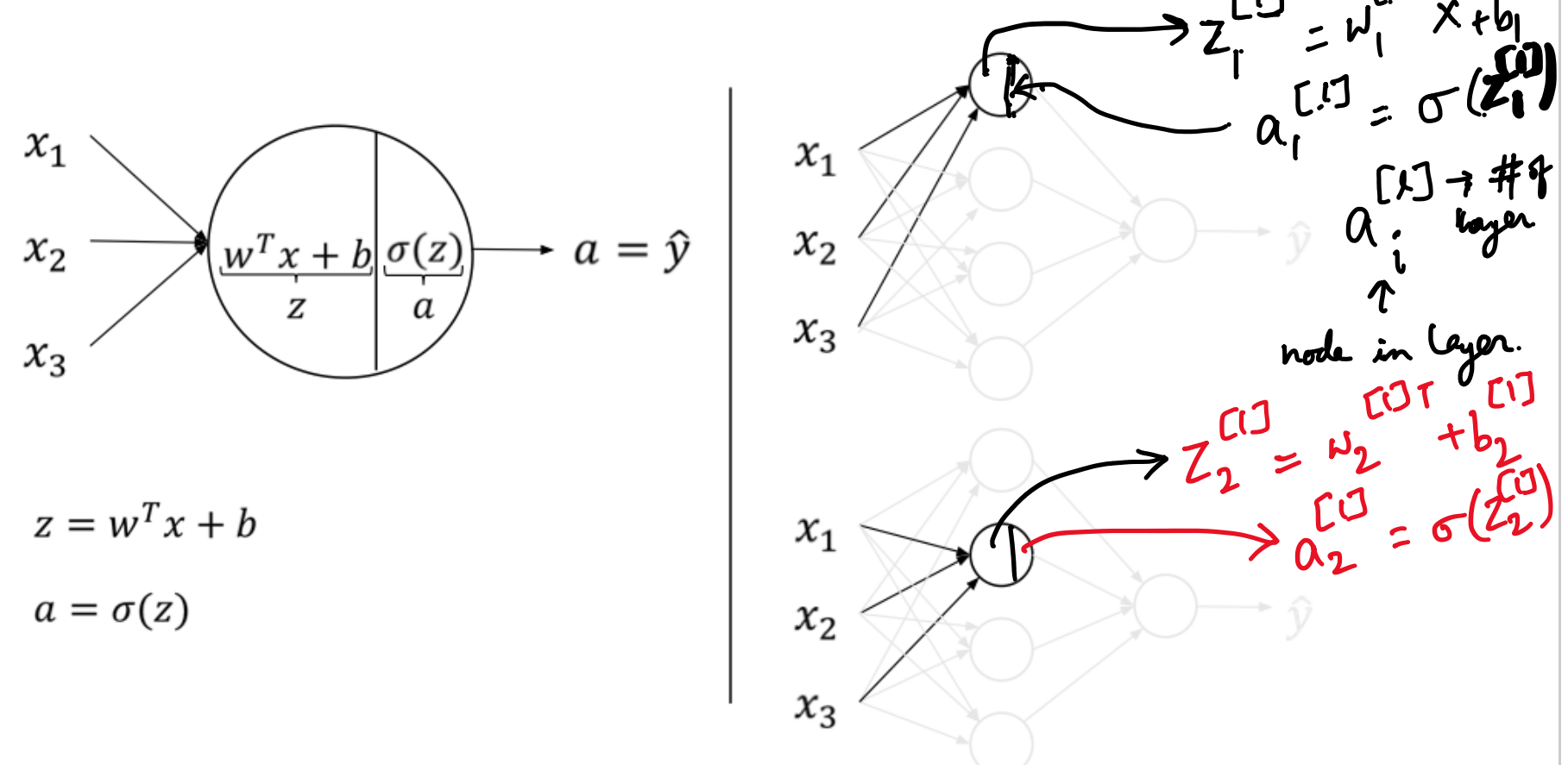
Neural Networks Representation



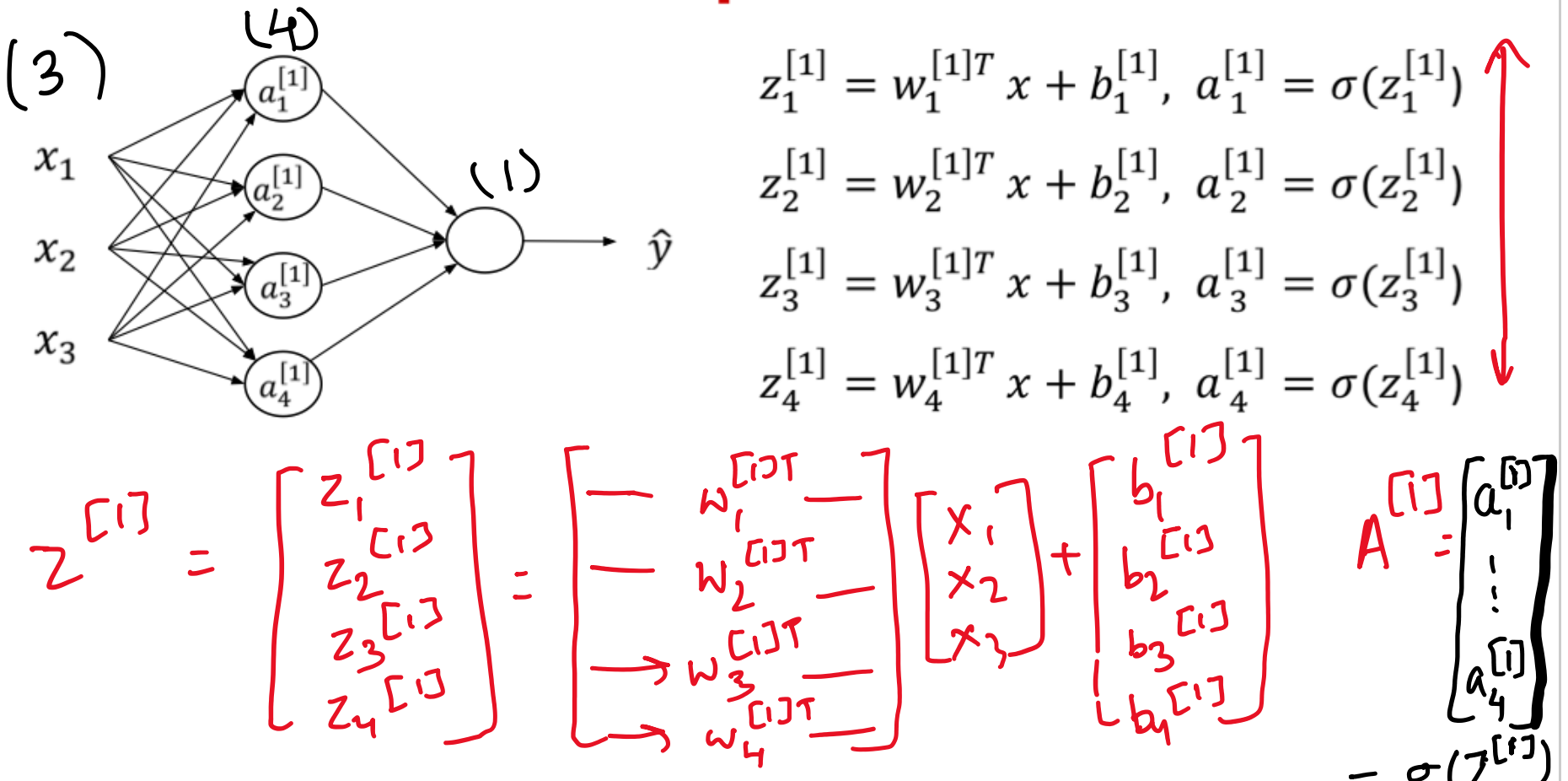
Neural Networks Representation



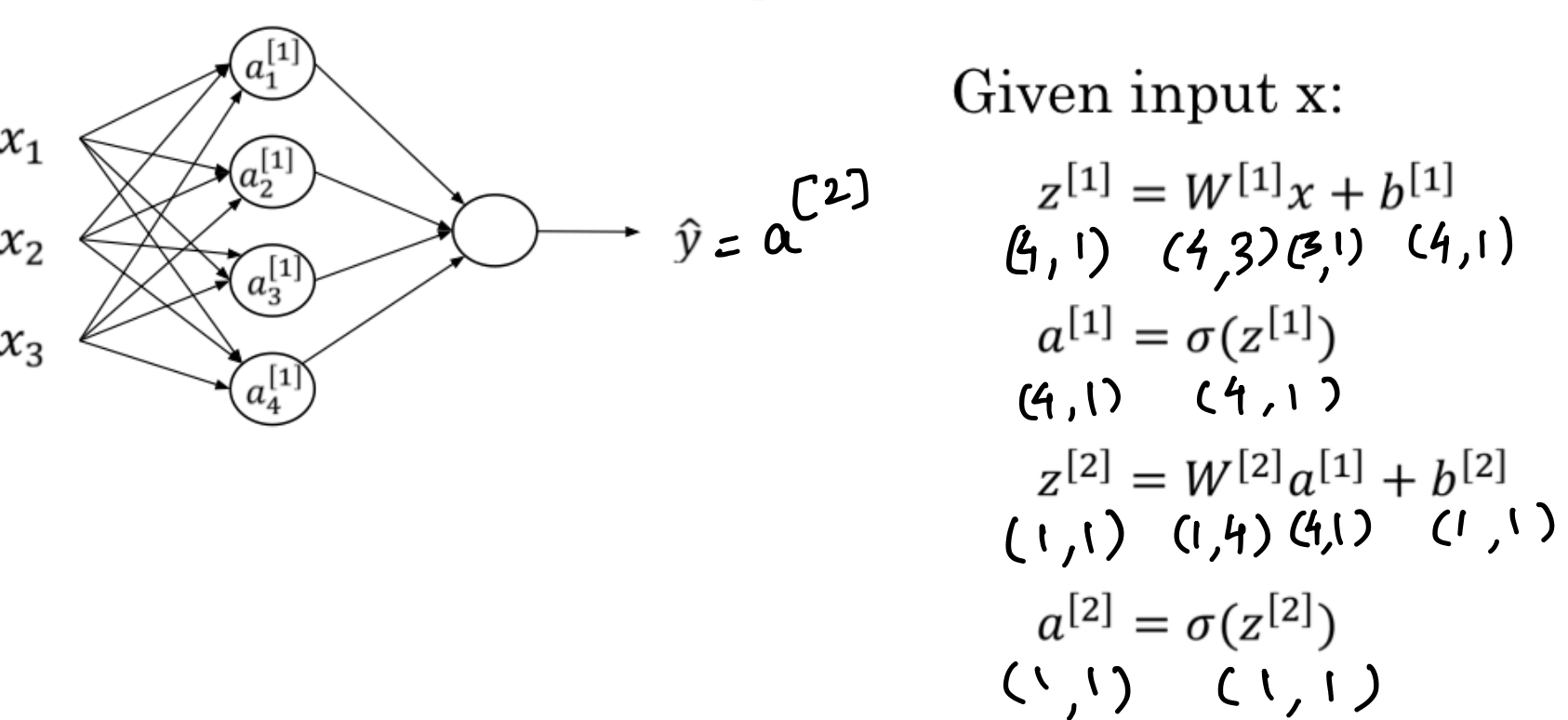
Neural Networks Representation



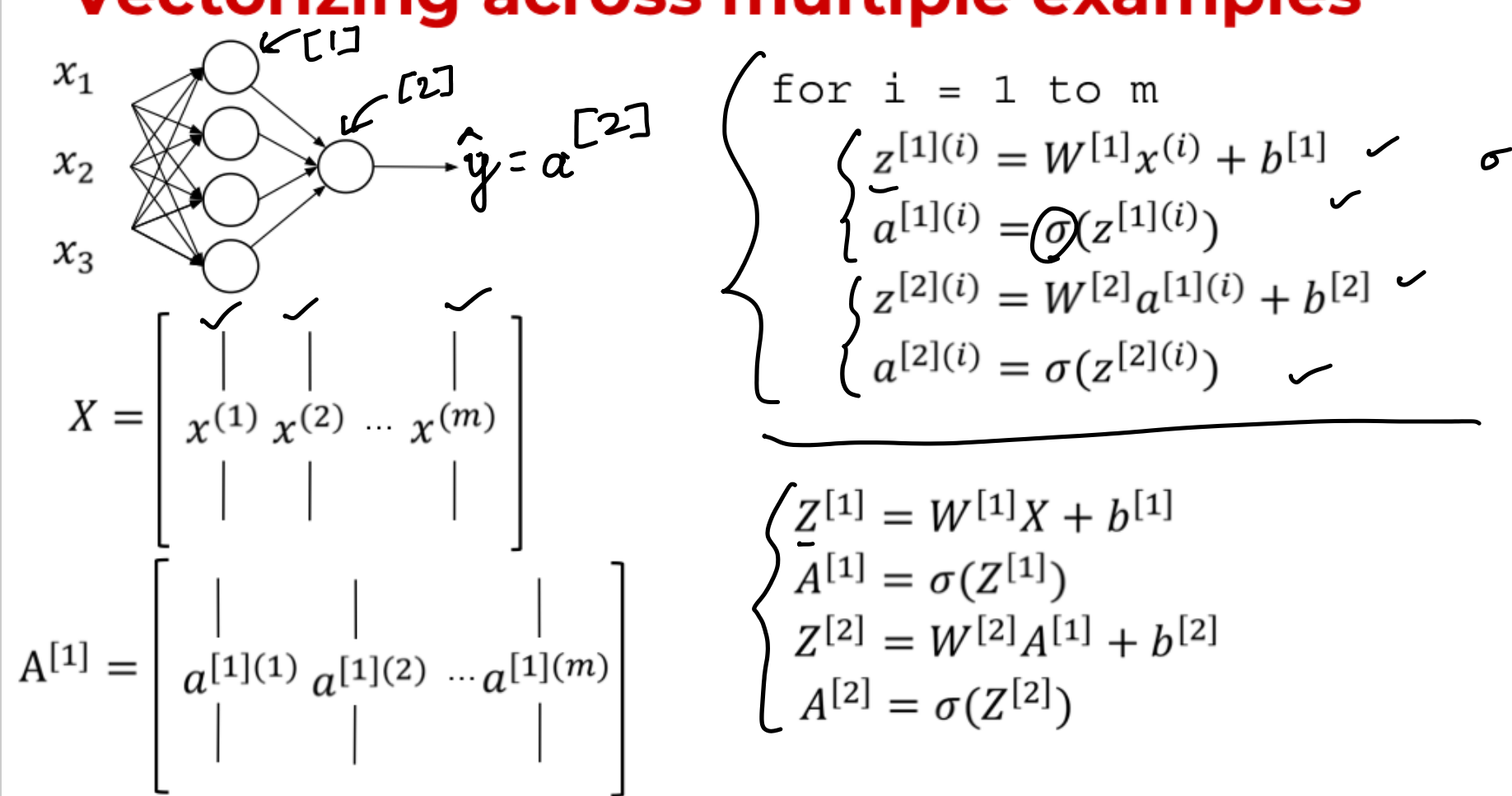
Neural Networks Representation



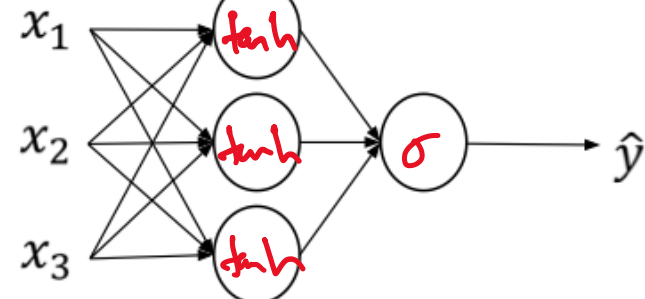
Neural Networks Representation Learning



Vectorizing across multiple examples



Activation functions



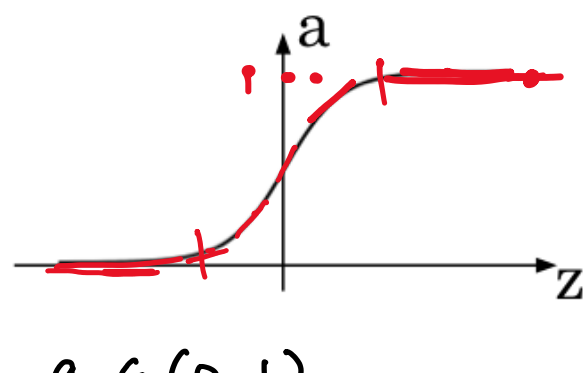
$$a = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Given x :

$$\begin{cases} z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[1]} = g^{[1]}(z^{[1]}) = \sigma(z^{[1]}) \\ z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} = g^{[2]}(z^{[2]}) = \sigma(z^{[2]}) \end{cases}$$

Sigmoid Activation Function

$W: W - \alpha \frac{dW}{dz}$

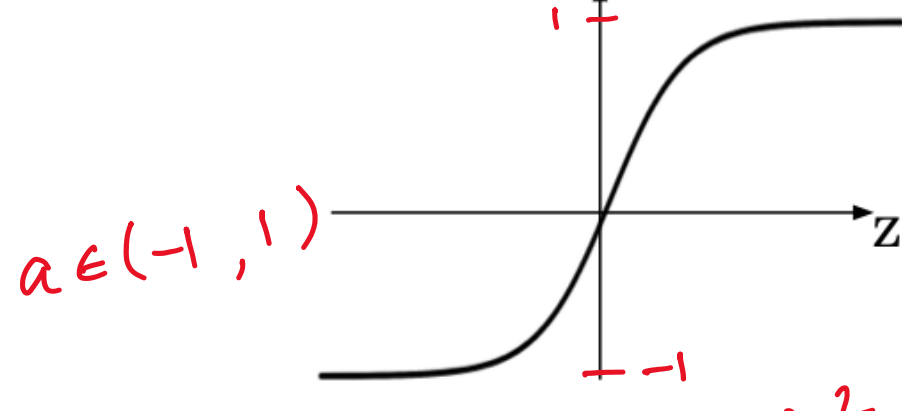


$$g(z) = \frac{1}{1 + e^{-z}}$$

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

for large z , $g'(z) \approx 0$
for large negative z , $g'(z) \approx 0$

Tanh Activation Function



$$a = g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - (g(z))^2$$

Tanh over Sigmoid

* Centralizes the data

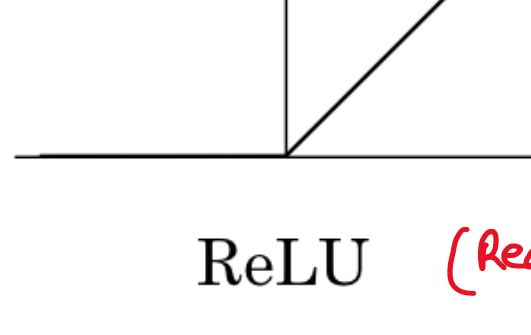
* Next layer's computation is easier

for $\uparrow z \Rightarrow g'(z) \approx 0$

for $\downarrow z \Rightarrow g'(z) \approx 0$

"Vanishing Gradient" Problem

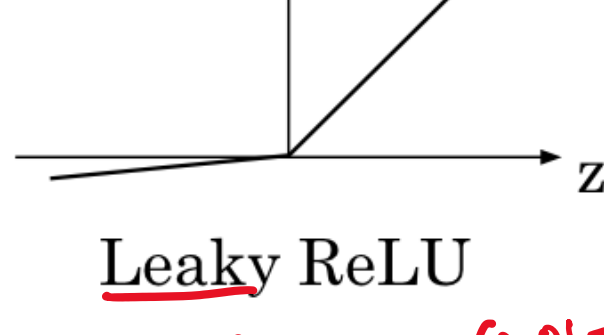
ReLU & Leaky ReLU Activation Functions



ReLU (Rectified Linear Unit)

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 0, & z < 0 \\ 1, & z > 0 \\ \text{undefined}, & z = 0 \end{cases}$$



Leaky ReLU

$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01, & z < 0 \\ 1, & z > 0 \end{cases}$$

* Tackles the vanishing grad. problem

Why non-linear activation functions?



Given x :

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

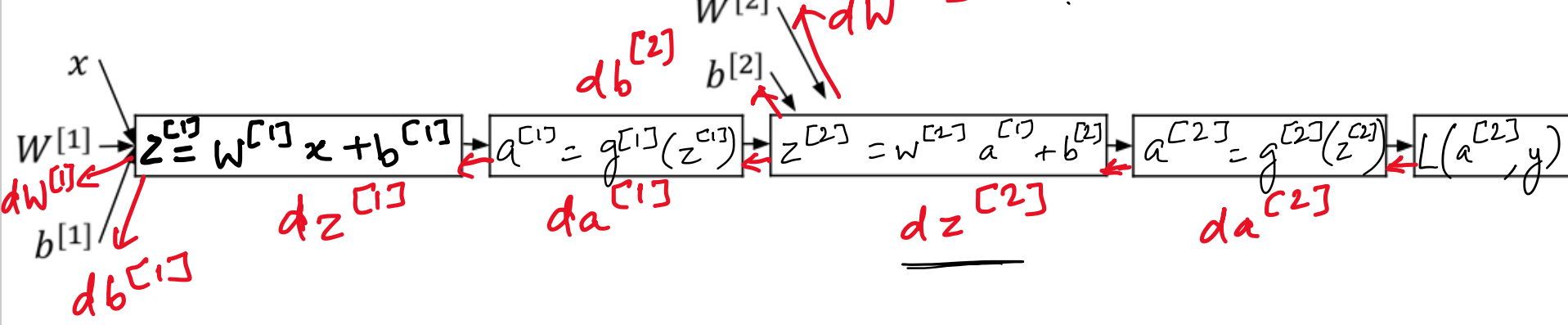
$$a^{[1]} = g^{[1]}(z^{[1]})$$

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

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

$$\begin{aligned} a^{[1]} &= z^{[1]} = W^{[1]}x + b^{[1]} \\ y &= a^{[2]} = z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ &= W^{[2]}(W^{[1]}x + b^{[1]}) + b^{[2]} \\ &= \underbrace{W^{[2]}W^{[1]}}_{W'}x + \underbrace{W^{[2]}b^{[1]} + b^{[2]}}_{b'} \\ &= W'x + b' \end{aligned}$$

Computation Graph



Gradient descent for neural networks

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

$$\text{Cost Function: } J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(a^{[2]}, y)$$

Gradient Descent:

Repeat {

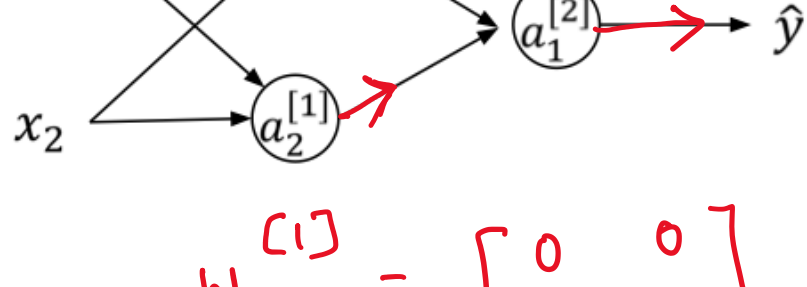
$$\begin{aligned} &\text{Compute } \frac{dW^{[1]}}{dz^{[1]}}, \frac{db^{[1]}}{dz^{[1]}}, \frac{dW^{[2]}}{dz^{[2]}}, \frac{db^{[2]}}{dz^{[2]}} \\ &W^{[1]}: \frac{dW^{[1]}}{dz^{[1]}} - \alpha \frac{dW^{[1]}}{dz^{[1]}}; \quad b^{[1]}: \frac{db^{[1]}}{dz^{[1]}} - \alpha \frac{db^{[1]}}{dz^{[1]}} \\ &W^{[2]}: \frac{dW^{[2]}}{dz^{[2]}} - \alpha \frac{dW^{[2]}}{dz^{[2]}}; \quad b^{[2]}: \frac{db^{[2]}}{dz^{[2]}} - \alpha \frac{db^{[2]}}{dz^{[2]}} \end{aligned}$$

Derivative formulas

Vectorized form

$$\begin{aligned} dz^{[2]} &= (a^{[2]} - y) & \rightarrow & dz^{[2]} = A^{[2]} - Y \\ dW^{[2]} &= dz^{[2]}a^{[1]T} & \rightarrow & dW^{[2]} = \frac{1}{m} dz^{[2]}A^{[1]T} \\ db^{[2]} &= dz^{[2]} & \rightarrow & db^{[2]} = \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis} = 1, \text{keepdims} = \text{True}) \\ dz^{[1]} &= W^{[2]T}dz^{[2]} * g^{[1]'}(z^{[1]}) & \rightarrow & dz^{[1]} = W^{[2]T}dz^{[2]} * g^{[1]'}(z^{[1]}) \\ dW^{[1]} &= dz^{[1]}x^T & \rightarrow & dW^{[1]} = \frac{1}{m} dz^{[1]}X^T \\ db^{[1]} &= dz^{[1]} & \rightarrow & db^{[1]} = \frac{1}{m} \text{np.sum}(dz^{[1]}, \text{axis} = 1, \text{keepdims} = \text{True}) \end{aligned}$$

Zero Initialization



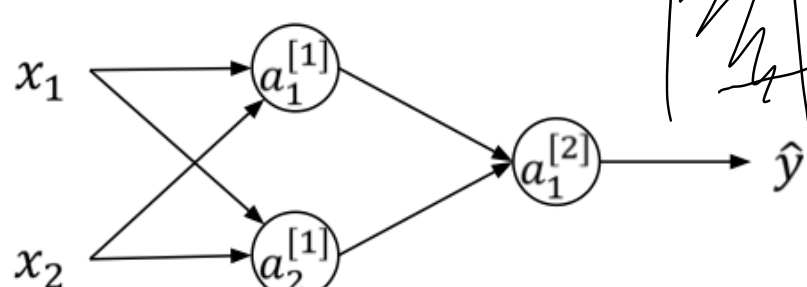
$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$a_1 = a_2$$

$$W^{[2]} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

* All hidden units/nodes will learn the same weights (Defeats the purpose of having multiple nodes)

Random Initialization



$$W^{[1]} = \text{np.random.randn}(2, 2) * 0.01$$

$$b^{[1]} = \text{np.zeros}(2, 1)$$

$$W^{[2]} = \text{np.random.randn}(1, 2) * 0.01$$

$$b^{[2]} = 0$$