

Deep learning

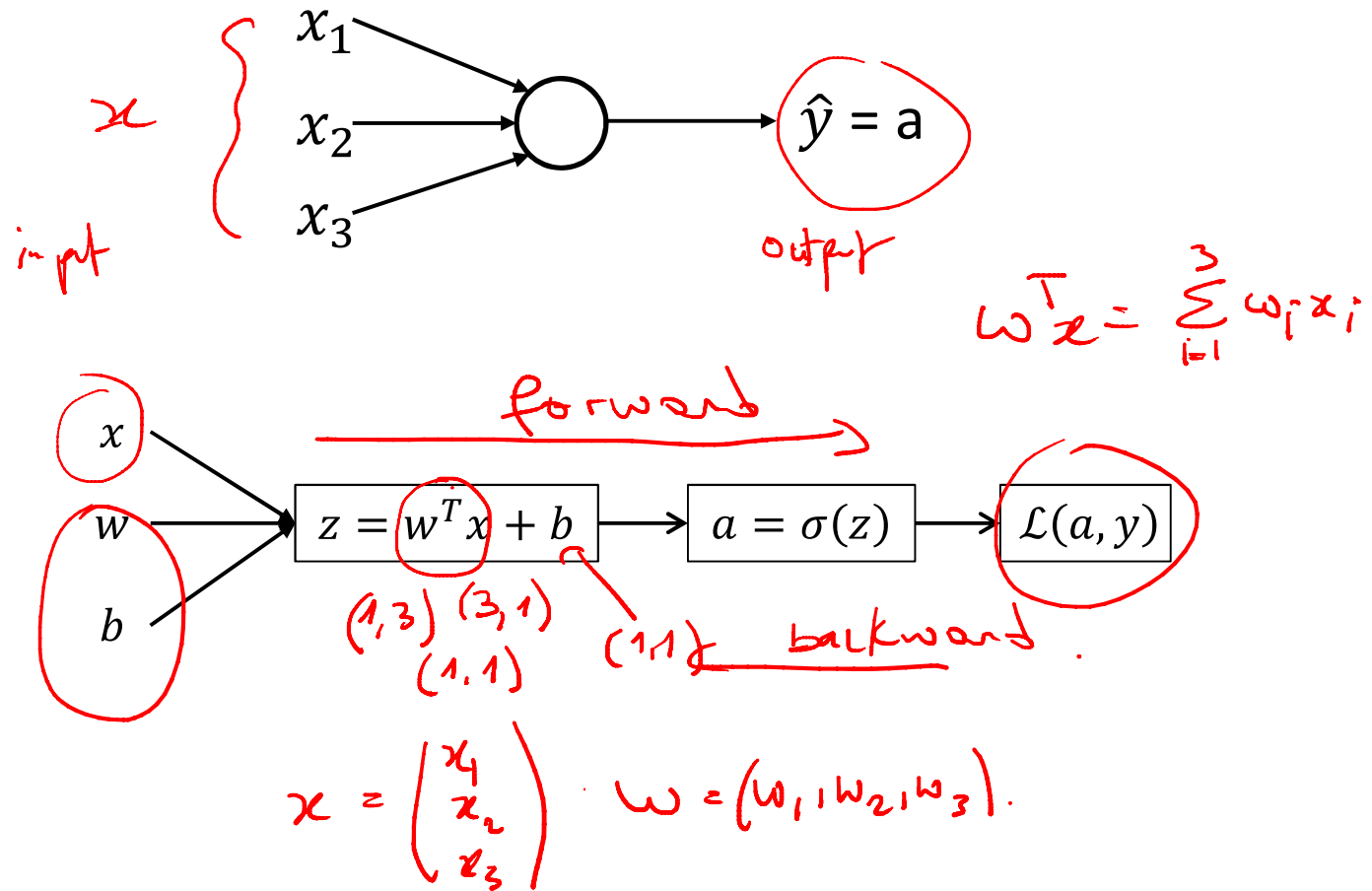
Dr. Aissa Boulmerka
a.boulmerka@centre-univ-mila.dz

2023-2024

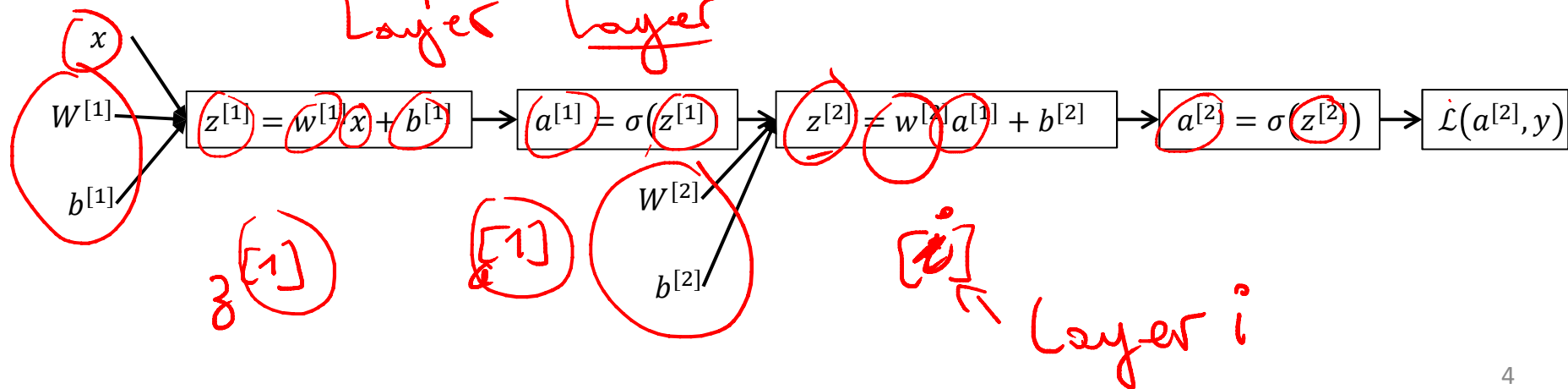
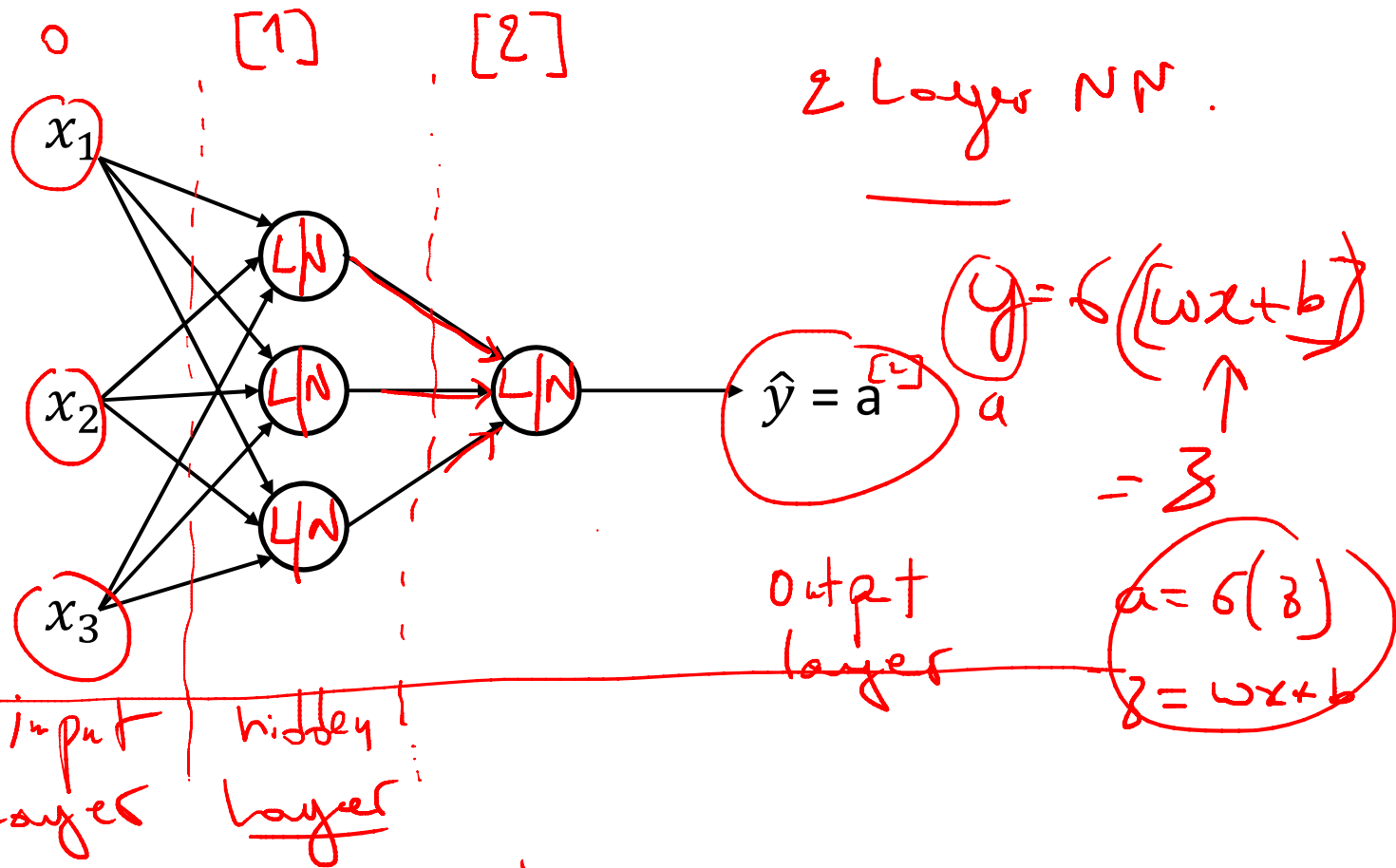
CHAPTER 2

SHALLOW NEURAL NETWORKS

What is a Neural Network?



What is a Neural Network?



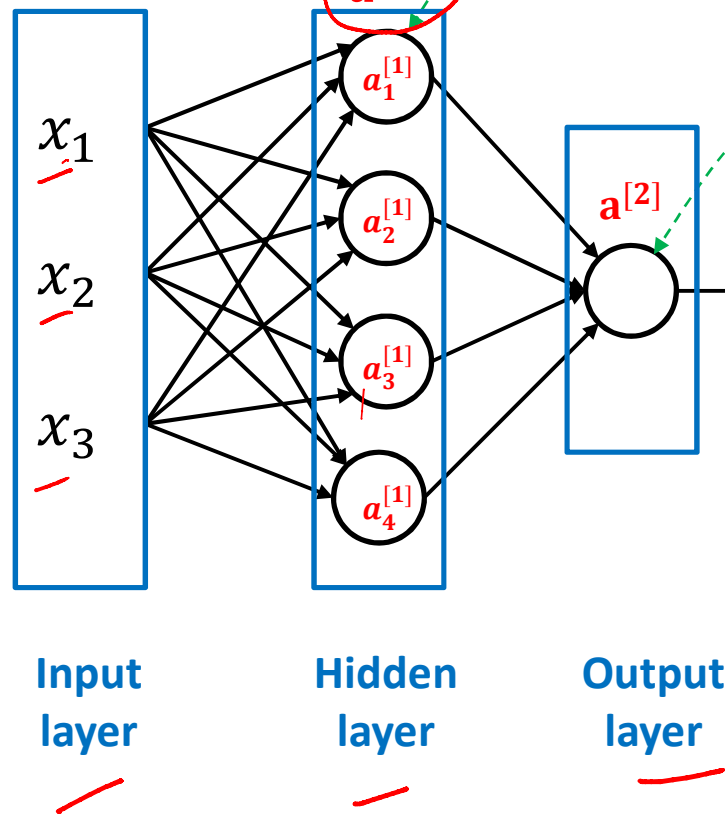
Neural Network Representation

2 layers neural network

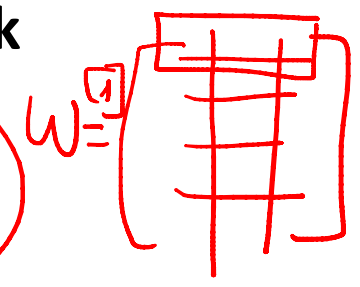
$$a^{[0]} = (a_1^{[0]}, a_2^{[0]}, a_3^{[0]}, a_4^{[0]})$$

$a^{[0]} = X$

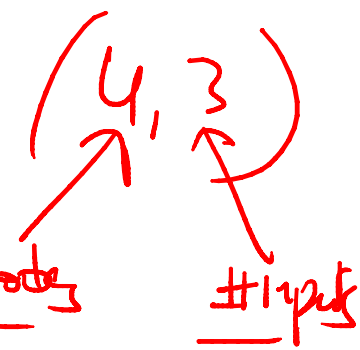
$$a_1^{[1]} = \sigma(z_1^{[1]})$$



$$w^{[1]} = (4, 3) \quad b^{[1]} = (4, 1)$$

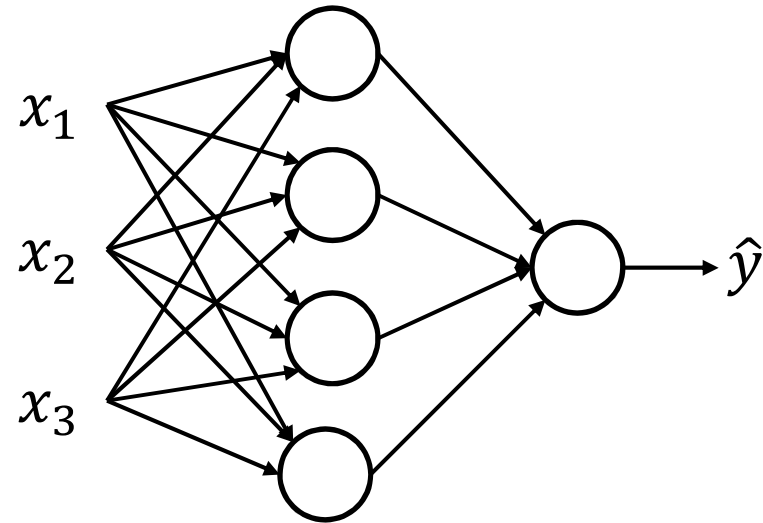
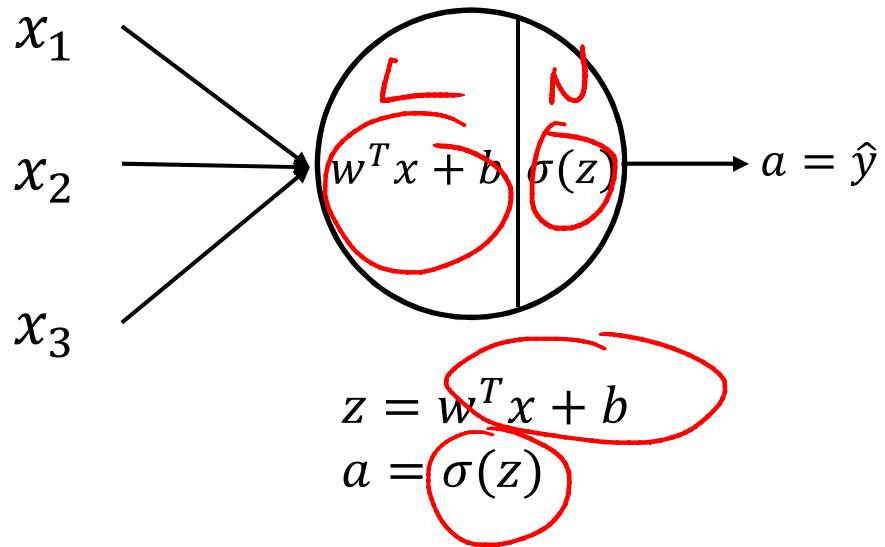


$$w^{[2]}, b^{[2]} = (1, 4) \quad (1, 1)$$

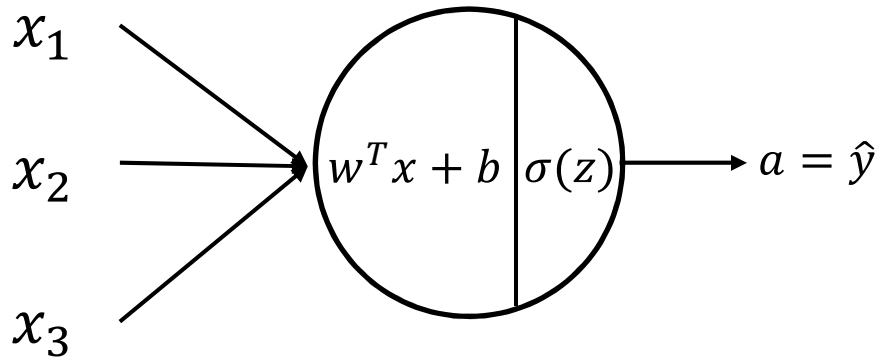


$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix}$$

Neural Network Representation



Neural Network Representation



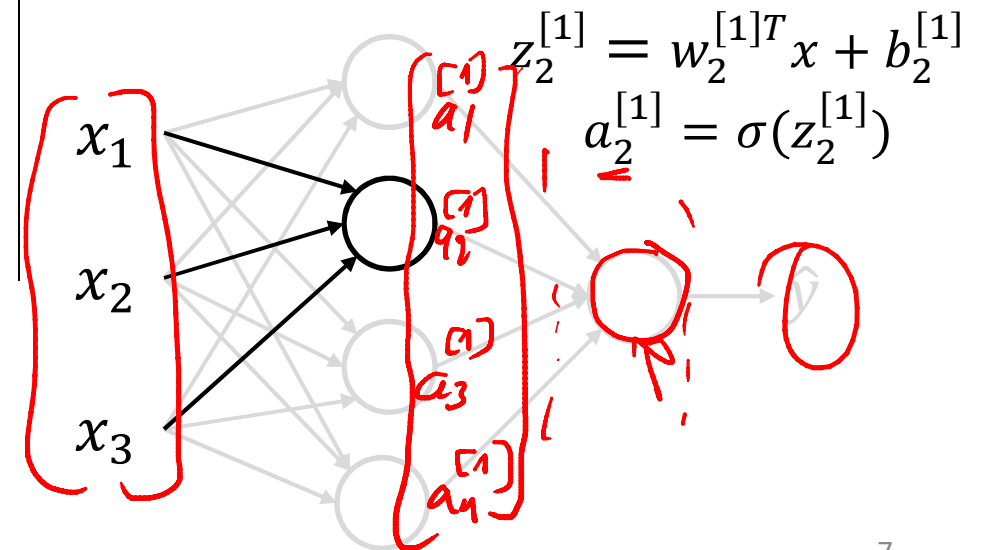
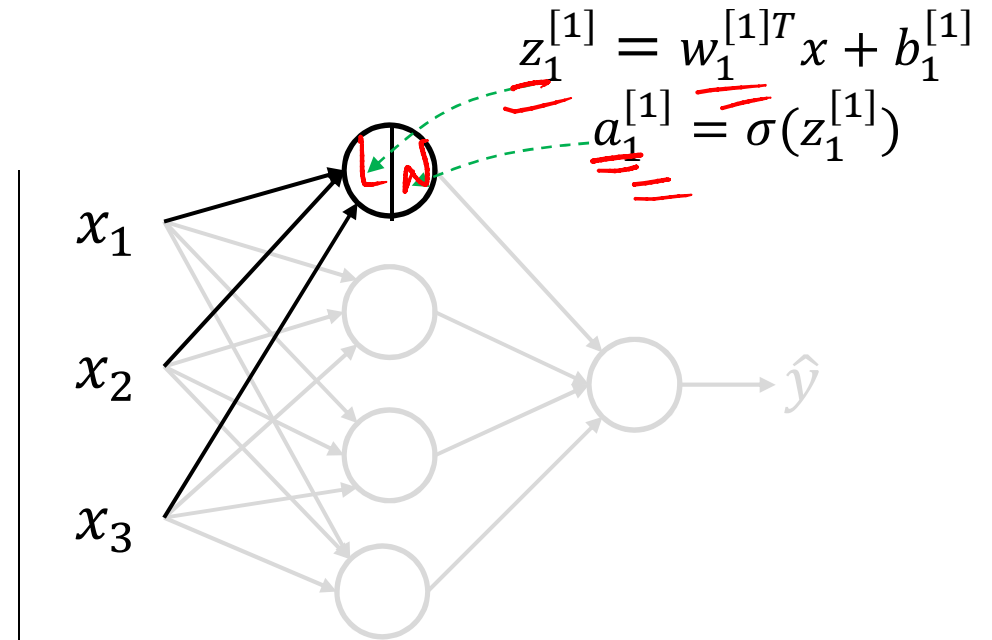
$$z = w^T x + b$$

$$a = \sigma(z)$$

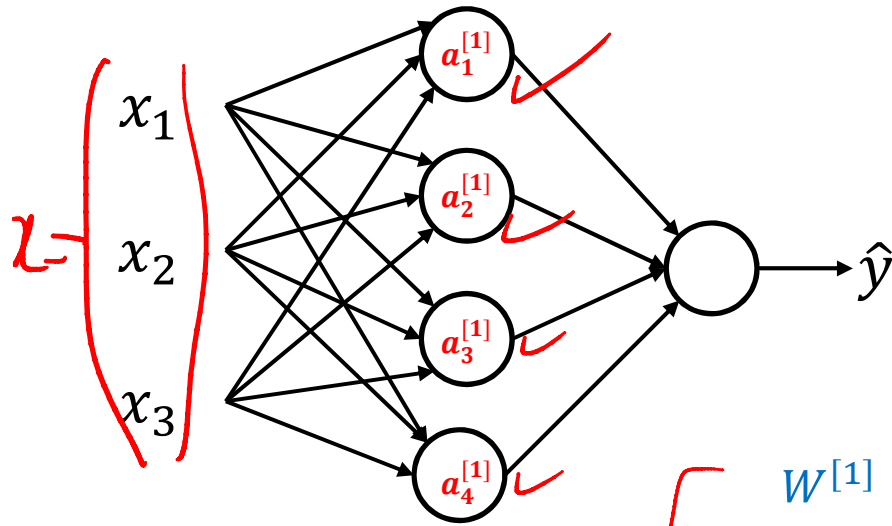
Handwritten notes in red:

$$z_1^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

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



Neural Network Representation



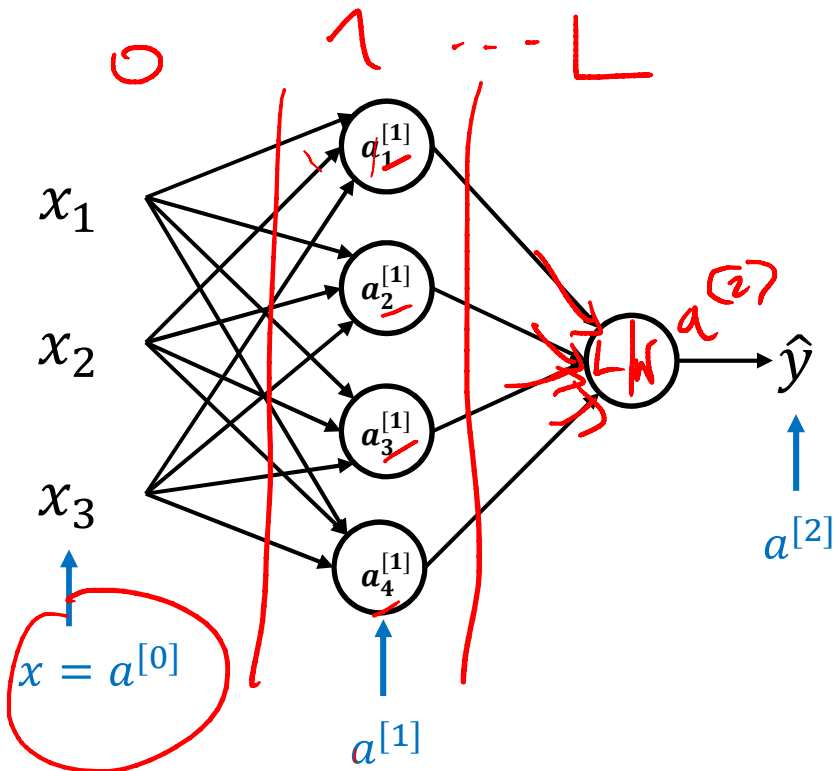
$$\begin{aligned} z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]}, & a_1^{[1]} &= \sigma(z_1^{[1]}) \\ z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]}, & a_2^{[1]} &= \sigma(z_2^{[1]}) \\ z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]}, & a_3^{[1]} &= \sigma(z_3^{[1]}) \\ z_4^{[1]} &= w_4^{[1]T} x + b_4^{[1]}, & a_4^{[1]} &= \sigma(z_4^{[1]}) \end{aligned}$$

$$z^{[1]} = \begin{matrix} (4,1) \\ \begin{bmatrix} - & w_1^{[1]T} & - \\ - & w_2^{[1]T} & - \\ - & w_3^{[1]T} & - \\ - & w_4^{[1]T} & - \end{bmatrix} \end{matrix} \begin{matrix} (3,1) \\ \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \end{matrix} + \begin{matrix} (4,1) \\ \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} \end{matrix} = \begin{matrix} \begin{bmatrix} w_1^{[1]T} x + b_1^{[1]} \\ w_2^{[1]T} x + b_2^{[1]} \\ w_3^{[1]T} x + b_3^{[1]} \\ w_4^{[1]T} x + b_4^{[1]} \end{bmatrix} \end{matrix} = \begin{matrix} \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix} \end{matrix}$$

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

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix} = \sigma(z^{[1]})$$

Neural Network Representation Learning



for $i=0$ to $L-1$

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

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

Given input x :

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

$(4,1) \quad (4,3)(3,1) \quad (4,1)$

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

$(4,1) \quad (4,1)$

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

$(1,1) \quad (1,4)(4,1) \quad (1,1)$

$$N \quad a^{[2]} = \sigma(z^{[2]}) = \hat{y}$$

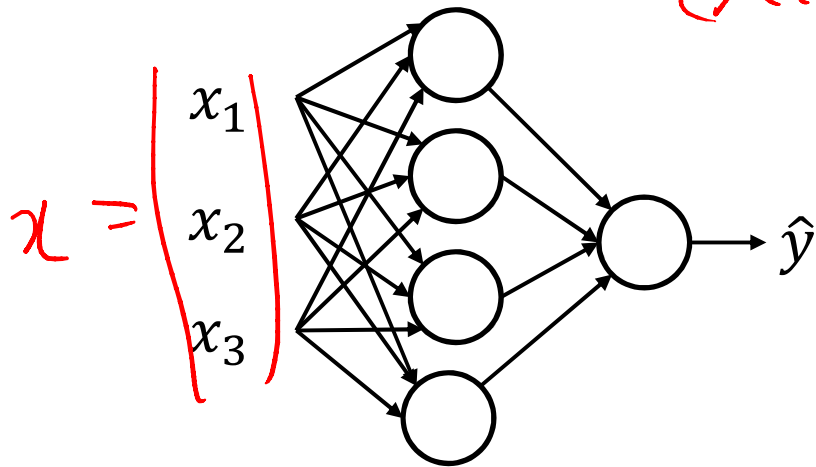
$(1,1) \quad (1,1)$

$x^{(1)}, x^{(2)}, \dots, x^{(m)}$

$(W, x_1^{(1)}, x_2^{(1)}, x_3^{(1)})$

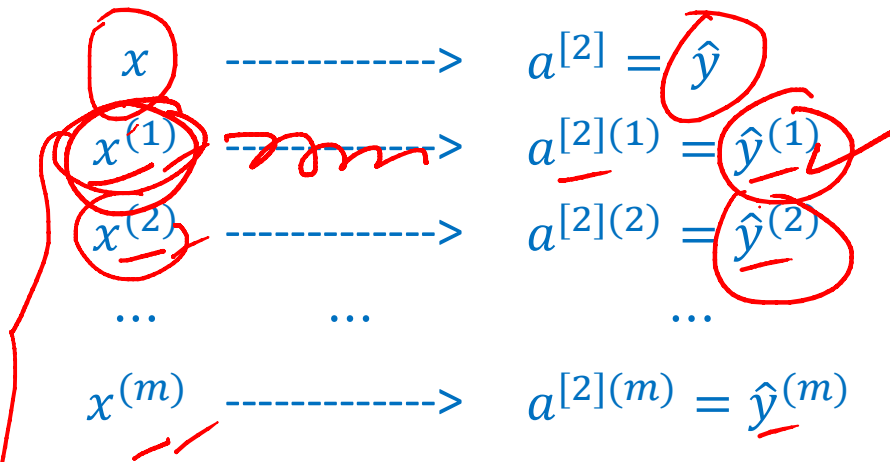
For loop across multiple examples

(X, Y) , $X = \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$



forward

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



for i=1 to m:

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

$\} [1](i)$

Layer 2 $\rightarrow a^{[2]}(i)$ \leftarrow Example (i)

Vectorizing across multiple examples

for i=1 to m:

$$z^{[1]}(i) = w^{[1]}x^{(i)} + b^{[1]}$$

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

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

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

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

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

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

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

$$X = \begin{bmatrix} | & | & \dots & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix}$$

(n_x, m)

$n_x = 3$

$$4 \begin{bmatrix} | & | & \dots & | \\ \dots & \dots & \dots & \dots \\ | & | & \dots & | \end{bmatrix}$$

$$Z^{[1]} = \begin{bmatrix} | & | & \dots & | \\ z^{[1]}(1) & z^{[1]}(2) & \dots & z^{[1]}(m) \\ | & | & \dots & | \end{bmatrix}$$

training examples (m)

$$A^{[1]} = \begin{bmatrix} | & | & \dots & | \\ a^{[1]}(1) & a^{[1]}(2) & \dots & a^{[1]}(m) \\ | & | & \dots & | \end{bmatrix}$$

hidden units

$m = 1000$

Justification for vectorized implementation

$$z^{1} = W^{[1]}x^{(1)} + b^{[1]} \quad , \quad z^{[1](2)} = W^{[1]}x^{(2)} + b^{[1]} \quad , \quad z^{[1](3)} = W^{[1]}x^{(3)} + b^{[1]}$$

$$W^{[1]}x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad \quad \quad W^{[1]}x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad \quad \quad W^{[1]}x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

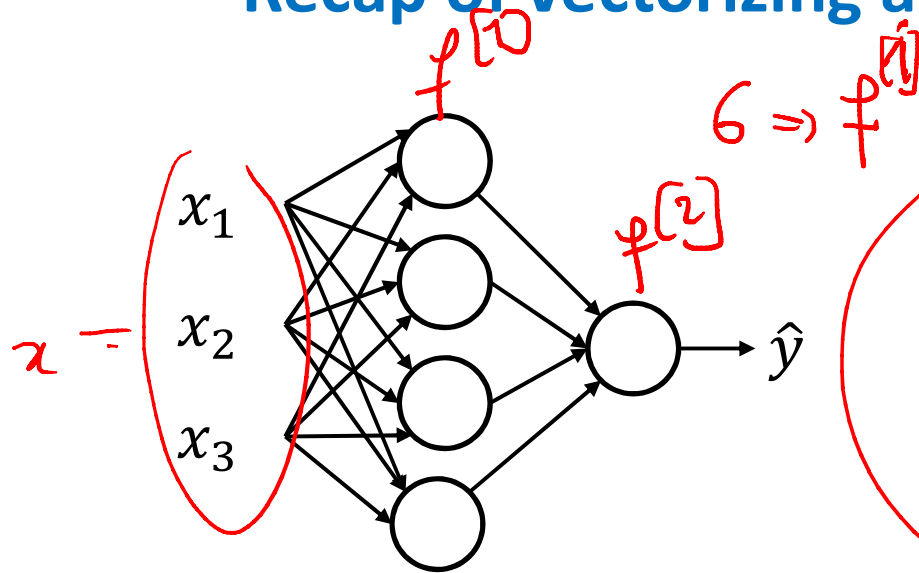
$$W^{[1]} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \\ | & | & | \end{bmatrix} + b^{[1]} = \begin{bmatrix} | & | & | \\ W^{[1]}x^{(1)} & W^{[1]}x^{(2)} & W^{[1]}x^{(3)} \\ | & | & | \\ + & + & + \\ b^{[1]} & b^{[1]} & b^{[1]} \end{bmatrix} = \begin{bmatrix} | & | & | \\ z^{1} & z^{[1](2)} & z^{[1](3)} \\ | & | & | \end{bmatrix}$$

X
 $Z^{[1]}$

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$m=3$

Recap of vectorizing across multiple examples



$$X = \begin{bmatrix} | & | & \dots & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} | & | & \dots & | \\ a^{[1]}(1) & a^{[1]}(2) & \dots & a^{[1]}(m) \\ | & | & \dots & | \end{bmatrix}$$

for $i=1$ to m :

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

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

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

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

activation

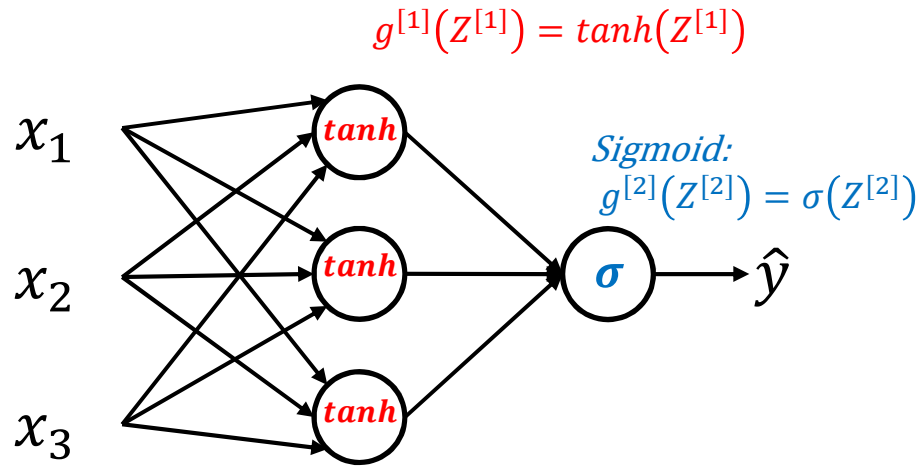
$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

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

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

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

Activation functions



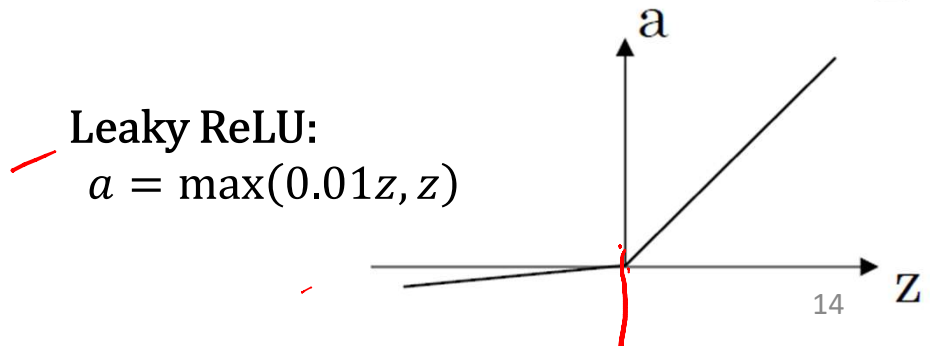
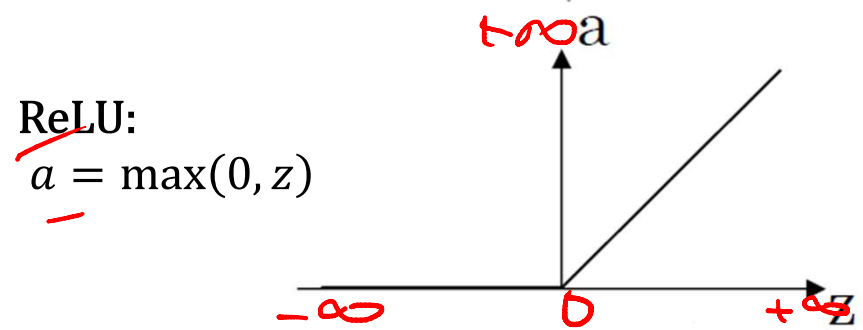
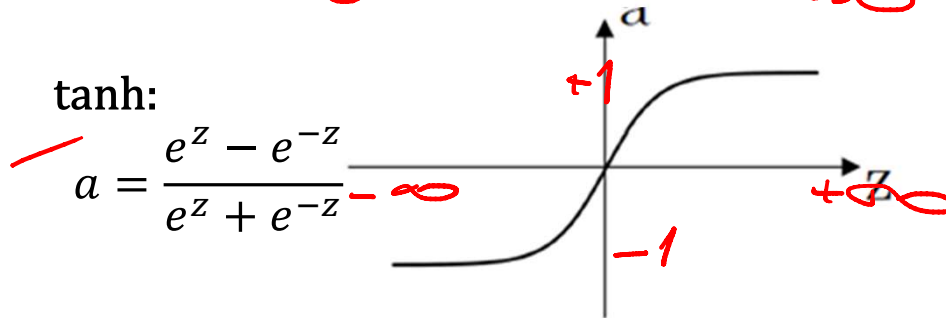
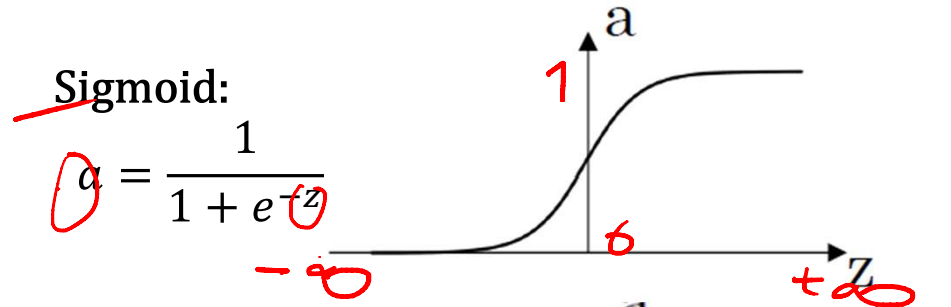
Given X:

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

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

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

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

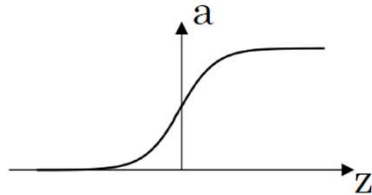


Pros and cons of activation functions

Sigmoid activation function:

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

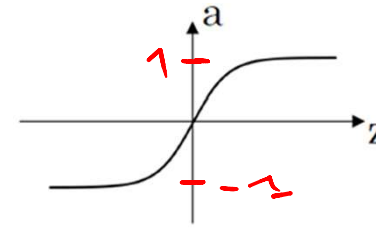
- Never use this, except for the output layer.
- if you are doing binary classification, or maybe almost never use this.



tanh activation function:

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

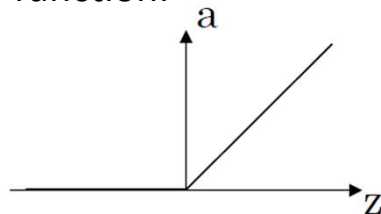
- The tanh is much strictly superior.



ReLU activation function:

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

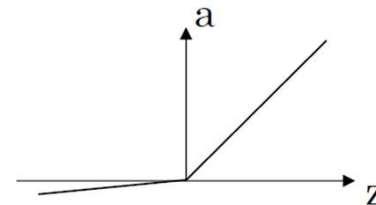
- The default and the most commonly used activation function is the ReLU.
- So if you're not sure what else to use, use the ReLU function.



Leaky ReLU activation function:

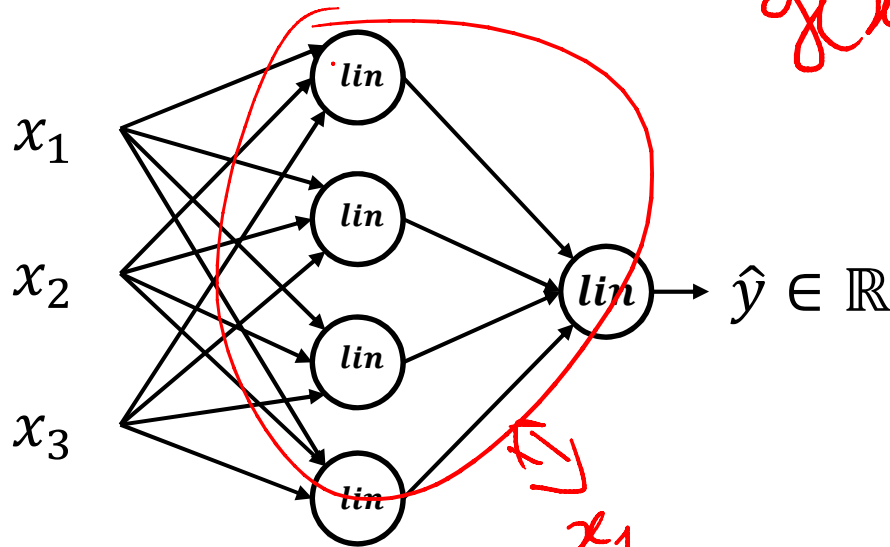
$$a = \max(0.01z, z)$$

- You can also try the leaky ReLU function.



Why do you need non linear activation functions?

$g(x) = \underline{x}$ $g(x) = \underline{ax + b}$



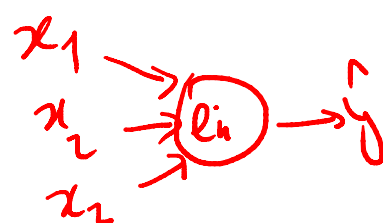
$$\underline{a}^{[1]} = \underline{z}^{[1]} = \underline{W}^{[1]}x + \underline{b}^{[1]}$$

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

$$\underline{a}^{[2]} = \underline{W}^{[2]}(\underline{W}^{[1]}x + \underline{b}^{[1]}) + \underline{b}^{[2]}$$

$$\underline{a}^{[2]} = \underline{(W^{[2]}W^{[1]})x + (W^{[2]}b^{[1]} + b^{[2]})}$$

$$\underline{a}^{[2]} = \underline{W'}x + \underline{b'}$$



Given X:

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

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

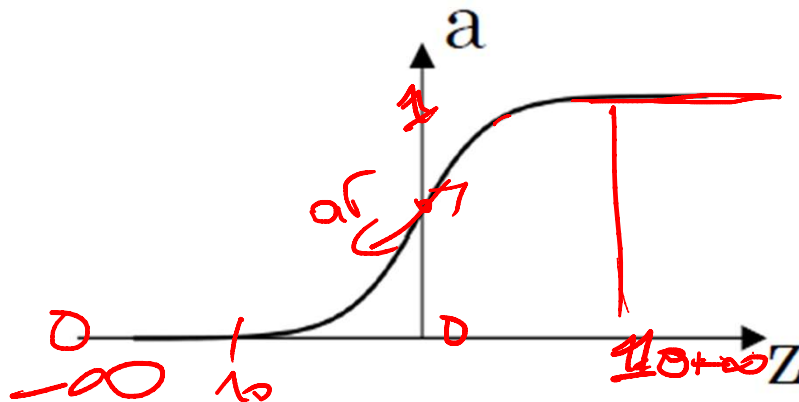
$$\underline{z}^{[2]} = \underline{W}^{[2]}\underline{a}^{[1]} + \underline{b}^{[2]}$$

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

Linear activation function:
 $g(z) = z$

Derivatives of activation functions

Sigmoid activation function:



$$g'(z) = \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z$$

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

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

$$a = g(z), \quad g'(z) = a(1 - a)$$

$$g(10) \approx 1$$

$$g'(10) \approx 0$$

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

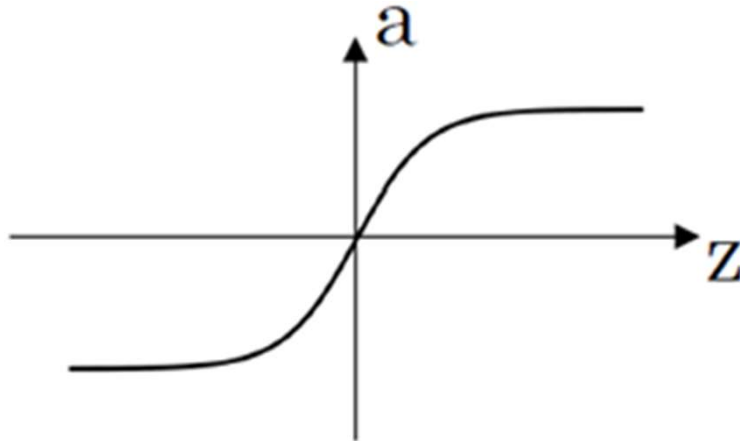
- $z = 10 \Rightarrow g(z) \approx 1$
- $\Rightarrow \frac{d}{dz} g(z) \approx 1(1 - 1) \approx 0$

- $z = -10 \Rightarrow g(z) \approx 0$
- $\Rightarrow \frac{d}{dz} g(z) \approx 0(1 - 0) \approx 0$

- $z = 0 \Rightarrow g(z) = 1/2$
- $\Rightarrow \frac{d}{dz} g(z) = \frac{1}{2} \left(1 - \frac{1}{2} \right) = 1/4$

Derivatives of activation functions

Tanh activation function:



$$g'(z) = \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z$$

$$= 1 - (\tanh(z))^2$$

$$a = g(z), \quad g'(z) = 1 - a^2$$

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

✓

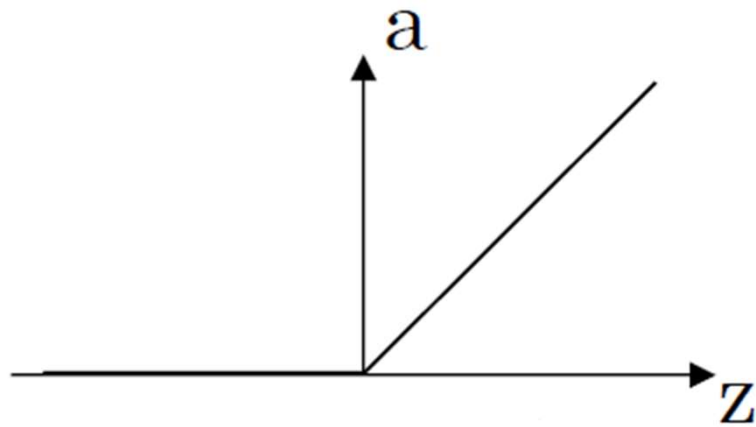
- $z = 10 \Rightarrow \tanh(z) \approx 1$
 $\Rightarrow g'(z) \approx 0$

- $z = -10 \Rightarrow g(z) \approx -1$
 $\Rightarrow g'(z) \approx 0$

- $z = 0 \Rightarrow g(z) = 0$
 $\Rightarrow g'(z) = 1$

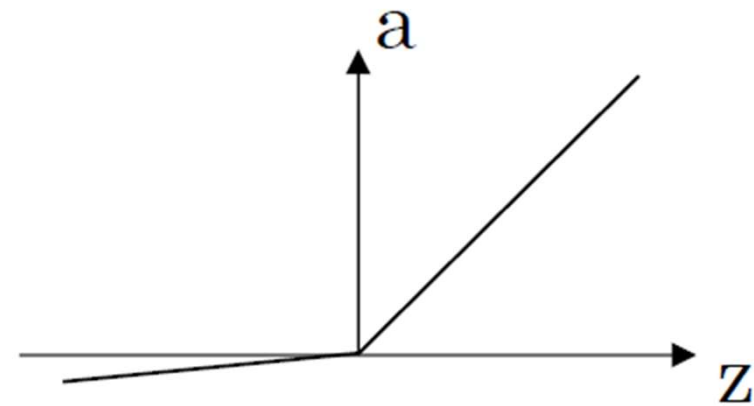
Derivatives of activation functions

ReLU and Leaky ReLU :



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

$$g'(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

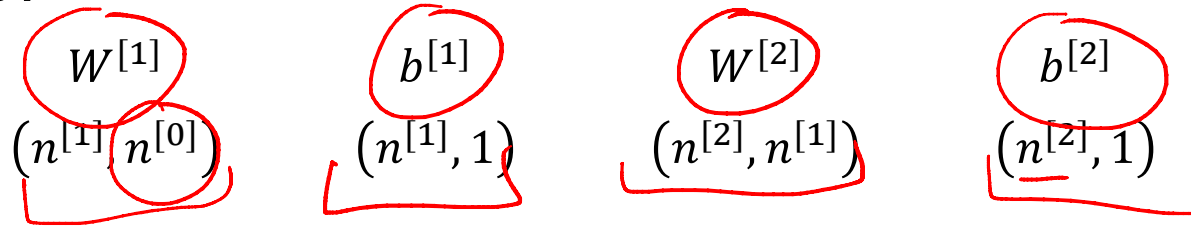


$$a = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0.01 & \text{if } z < 0 \end{cases}$$

Gradient descent for neural networks

Parameters :



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

Cost function:

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

output

Gradient descent:

Initialization of $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$?

```

Repeat {
  Compute predictions ( $\hat{y}^{(i)}, i = 1, \dots, m$ )
   $dW^{[1]} = \frac{\partial J}{\partial W^{[1]}}$ ,  $db^{[1]} = \frac{\partial J}{\partial b^{[1]}}$  derivatives
   $W^{[1]} = W^{[1]} - \alpha dW^{[1]}$ 
   $b^{[1]} = b^{[1]} - \alpha db^{[1]}$ 
  -----
   $dW^{[2]} = \frac{\partial J}{\partial W^{[2]}}$ ,  $db^{[2]} = \frac{\partial J}{\partial b^{[2]}}$  derivatives
   $W^{[2]} = W^{[2]} - \alpha dW^{[2]}$ 
   $b^{[2]} = b^{[2]} - \alpha db^{[2]}$ 
} Converge.
  
```

MAJ

Formulas for computing derivatives

Forward propagation:

$$\underline{Z^{[1]}} = \underline{W^{[1]}}X + \underline{b^{[1]}}$$

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

$$\underline{Z^{[2]}} = \underline{W^{[2]}}A^{[1]} + \underline{b^{[2]}}$$

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

Back propagation:

$\hat{y} - y$

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

$$L(g, y) = \frac{1}{2}(\hat{y} - y)^2$$

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

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

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

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

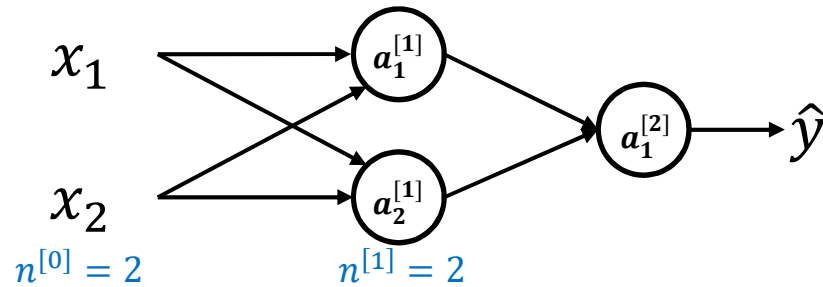
Element wise product

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

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

$(n^{[1]}, 1)$ $(n^{[1]},)$ reshape

What happens if you initialize weights to zero?



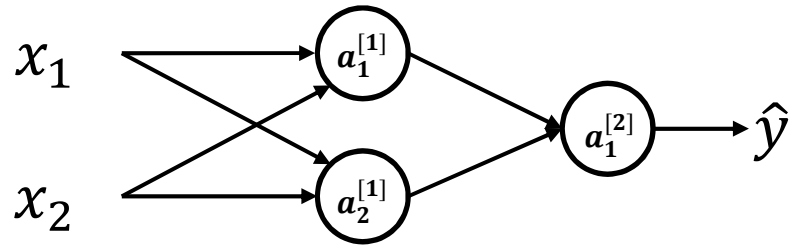
$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]} \text{ (symmetric)} \quad dz_1^{[1]} = dz_2^{[1]}$$

$$dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix} \quad W^{[1]} = W^{[1]} - \alpha dw$$

- The bias terms b can be initialized by 0, but initializing W to all 0s is a problem:
 - The two activations $a_1^{[1]}$ and $a_2^{[1]}$ will be the same, because both of these hidden units are computing exactly the same function.
 - After every single iteration of training the two hidden units are still computing exactly the same function.

Random initialization



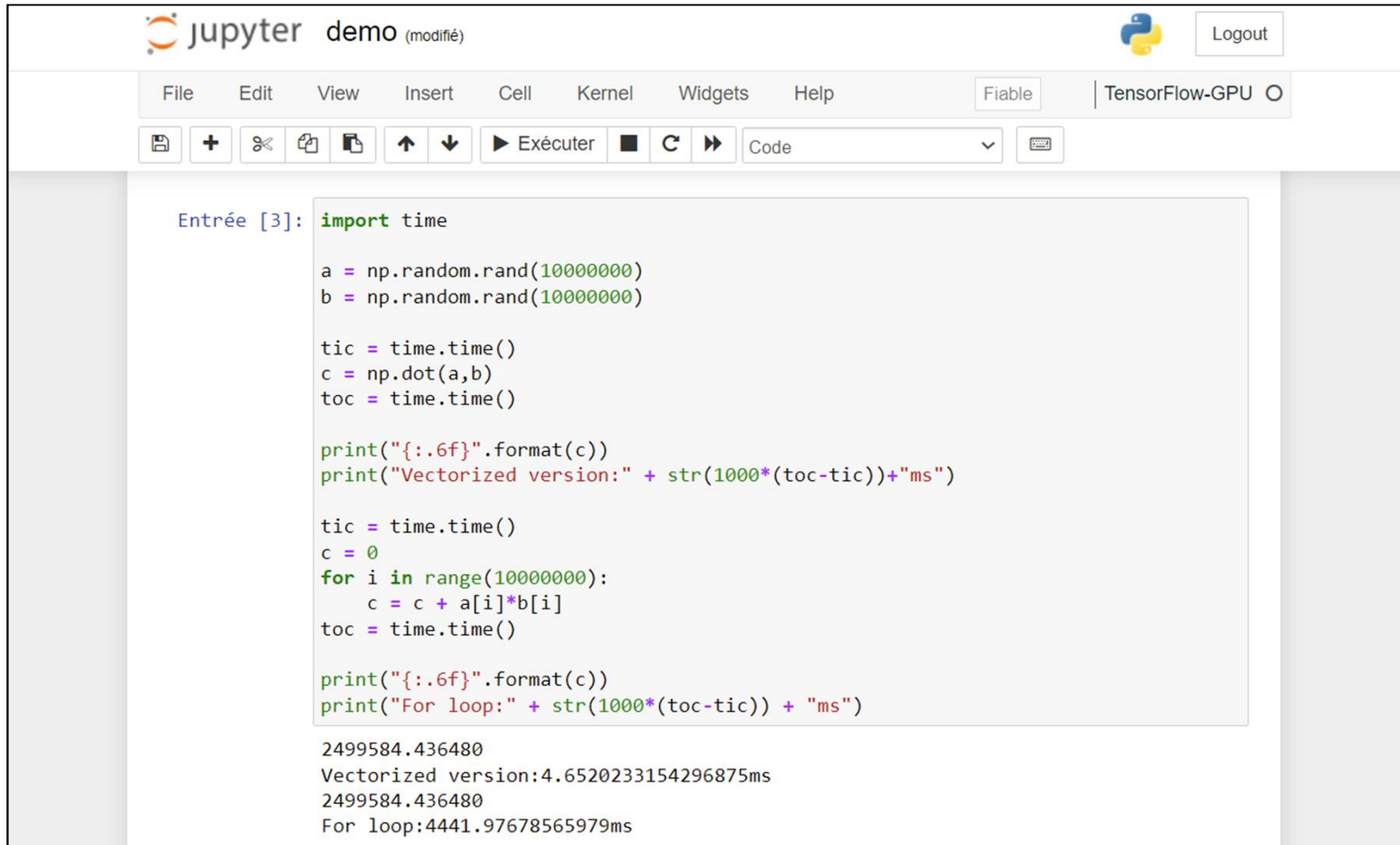
$$W^{[1]} = np.random.randn((2,2)) * 0.01$$

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

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

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

Vectorization demo



The screenshot shows a Jupyter Notebook interface with the following components:

- Header: "jupyter demo (modifié)" and a "Logout" button.
- Menu: File, Edit, View, Insert, Cell, Kernel, Widgets, Help.
- Toolbar: Includes icons for file operations, execution ("Exécuter"), and code editing.
- Code Cell: Labeled "Entrée [3]:", containing Python code that compares a vectorized dot product with a for-loop implementation.
- Output: Shows the results of the code execution, including timing information for both methods.

```
Entrée [3]: import time

a = np.random.rand(10000000)
b = np.random.rand(10000000)

tic = time.time()
c = np.dot(a,b)
toc = time.time()

print("{:.6f}".format(c))
print("Vectorized version:" + str(1000*(toc-tic))+ "ms")

tic = time.time()
c = 0
for i in range(10000000):
    c = c + a[i]*b[i]
toc = time.time()

print("{:.6f}".format(c))
print("For loop:" + str(1000*(toc-tic)) + "ms")

2499584.436480
Vectorized version:4.6520233154296875ms
2499584.436480
For loop:4441.97678565979ms
```


References

- Andrew Ng. Deep learning. Coursera.
- Geoffrey Hinton. Neural Networks for Machine Learning.
- Kevin P. Murphy. Probabilistic Machine Learning An Introduction. MIT Press, 2022.