

Deep learning

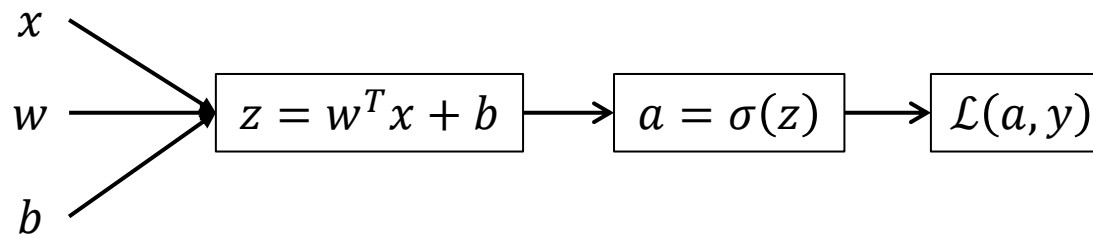
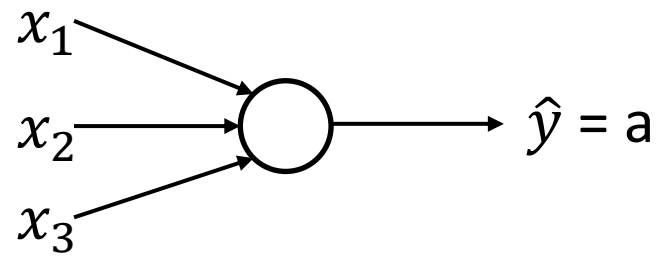
Dr. Aissa Boulmerka
a.boulmerka@centre-univ-mila.dz

2023-2024

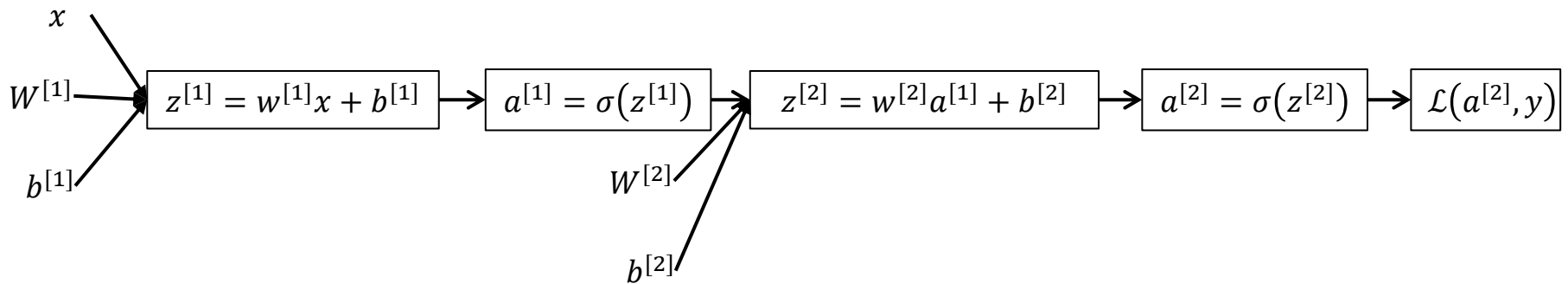
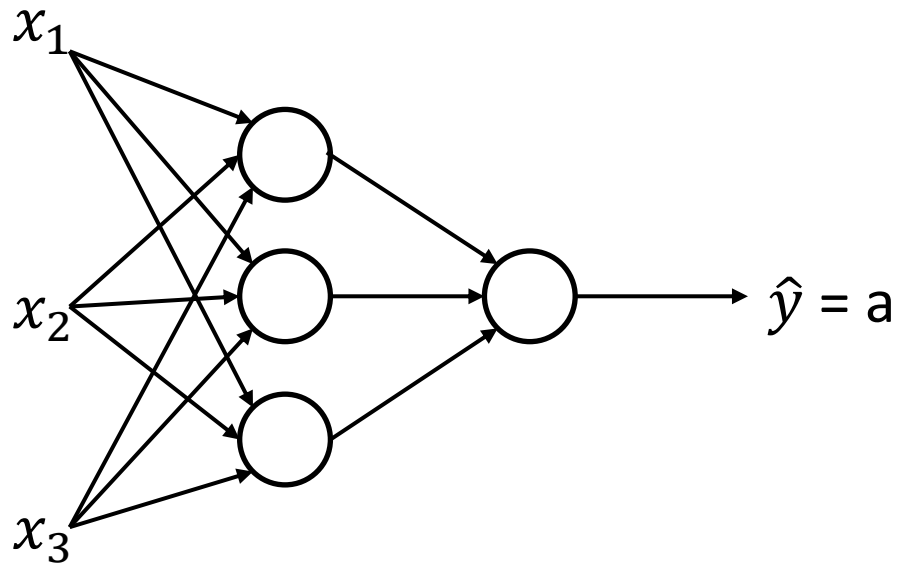
CHAPTER 2

SHALLOW NEURAL NETWORKS

What is a Neural Network?

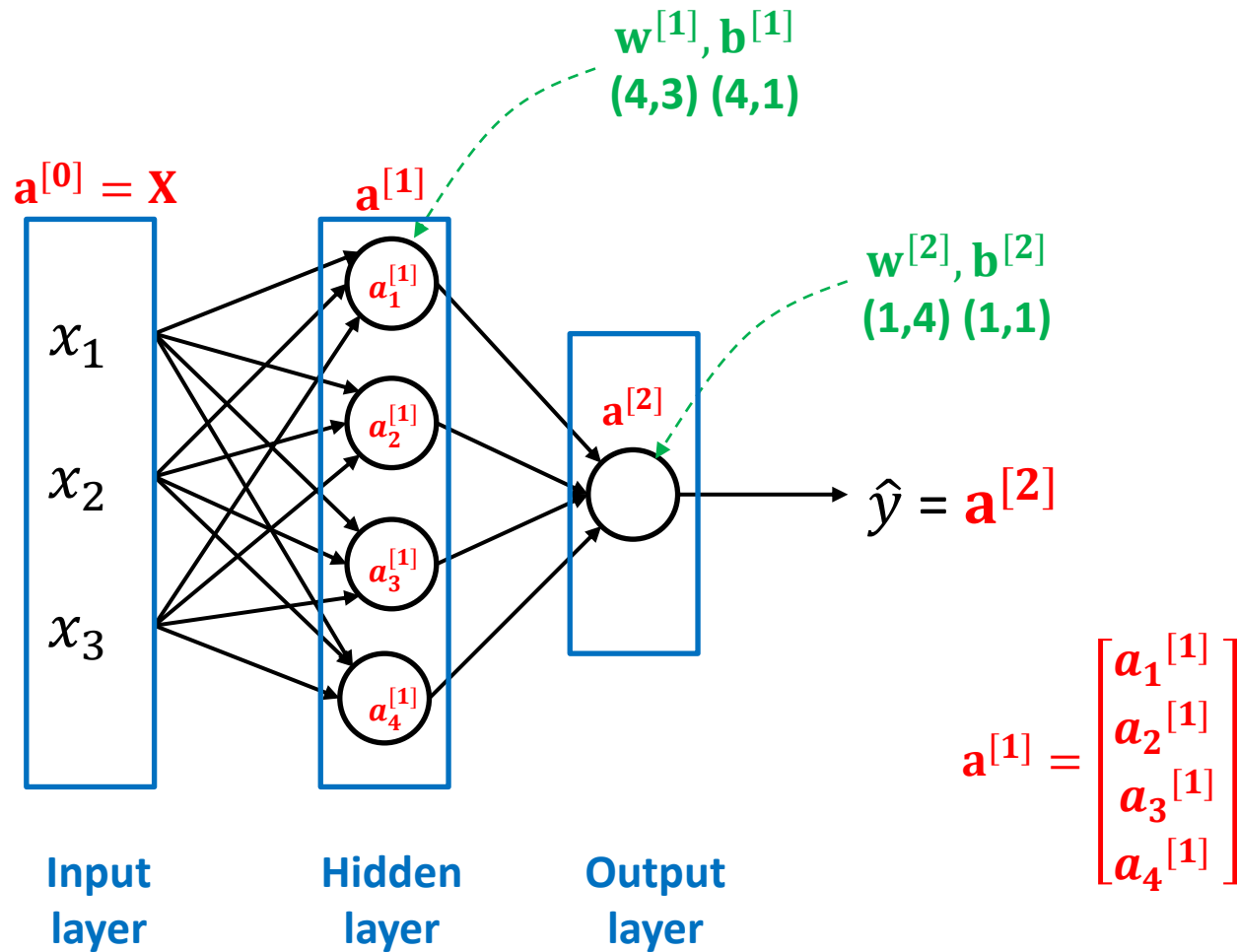


What is a Neural Network?

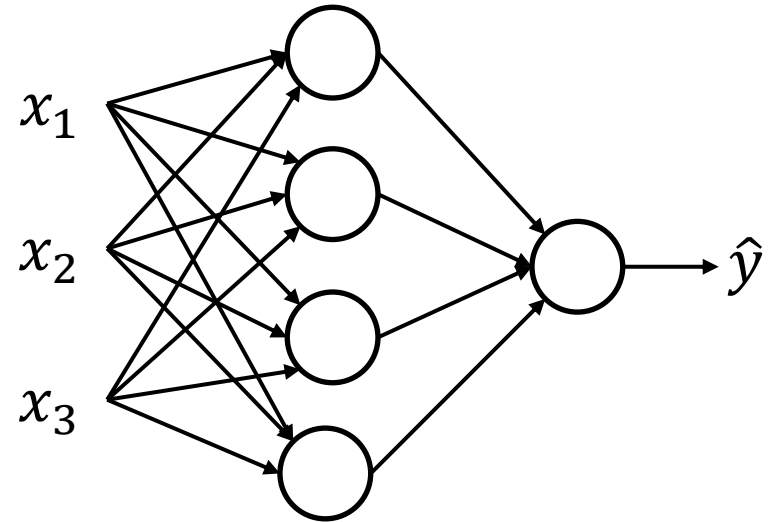
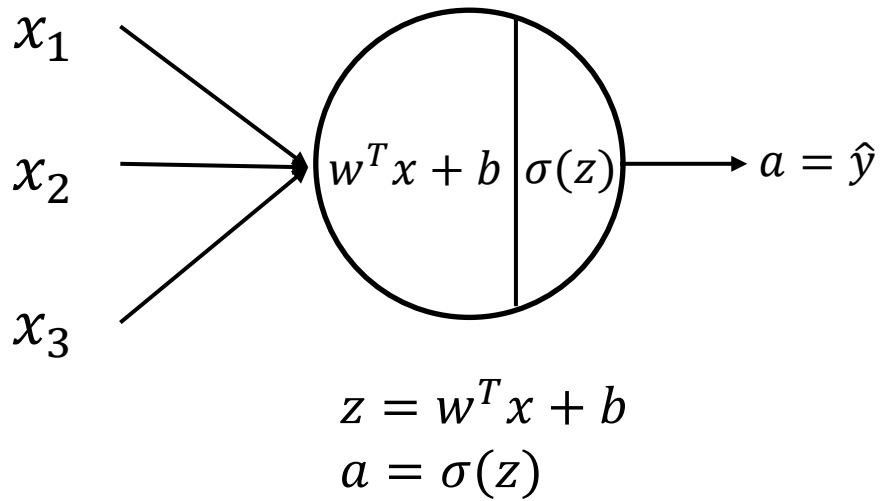


Neural Network Representation

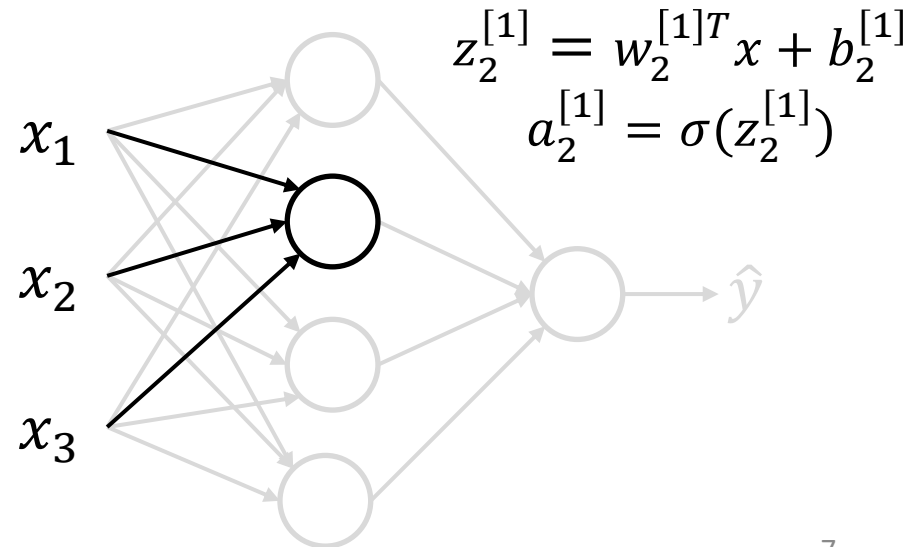
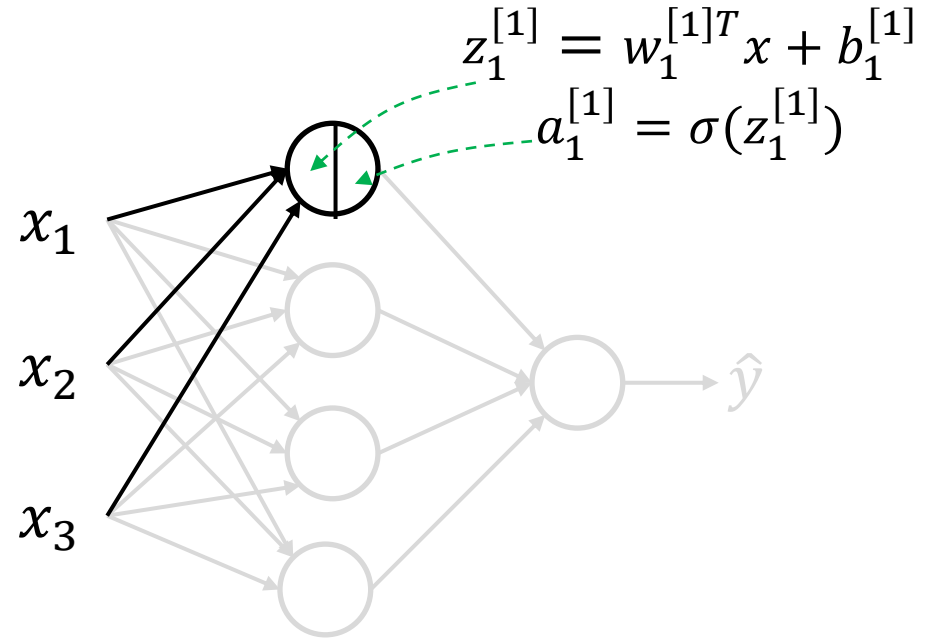
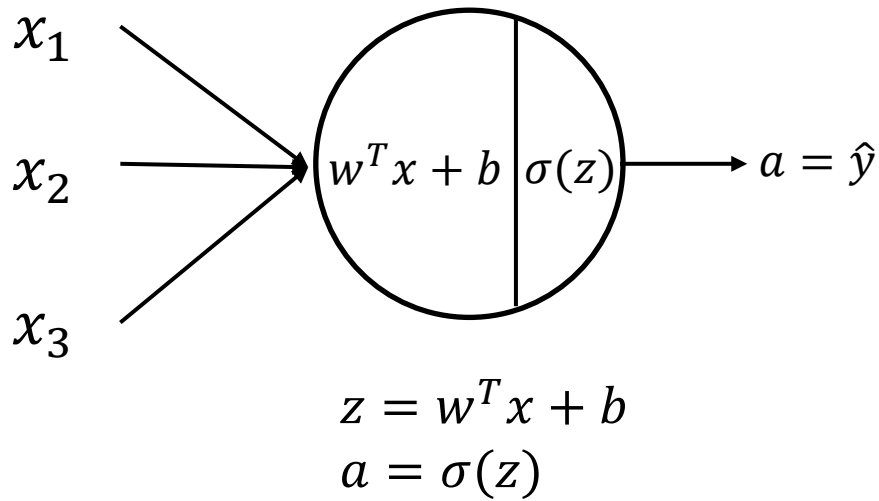
2 layers neural network



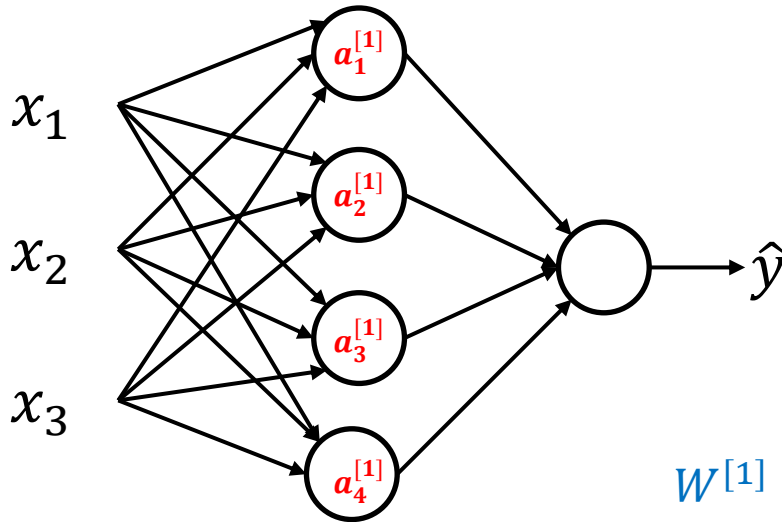
Neural Network Representation



Neural Network Representation



Neural Network Representation



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

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

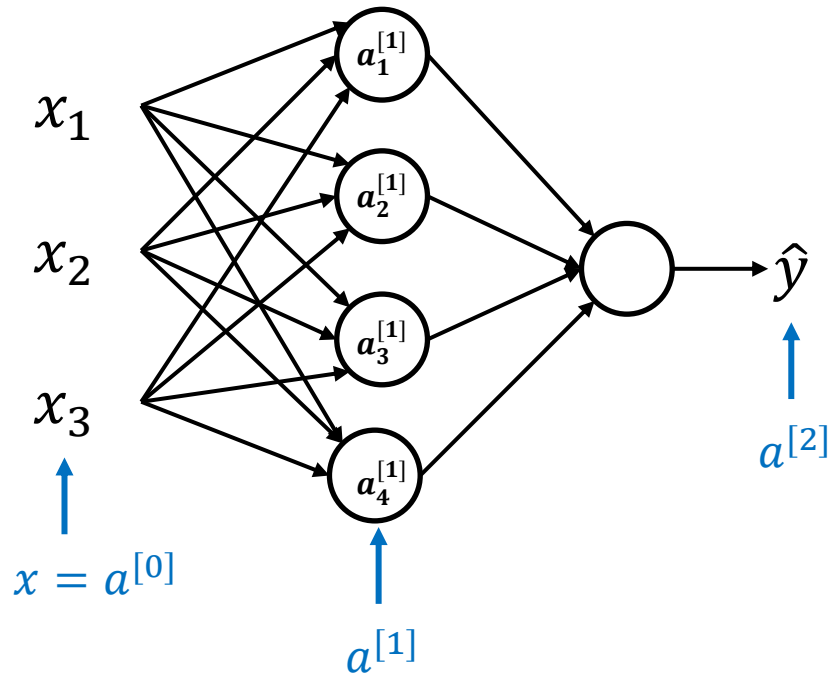
$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

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

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

Neural Network Representation Learning



Given input x :

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

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

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

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

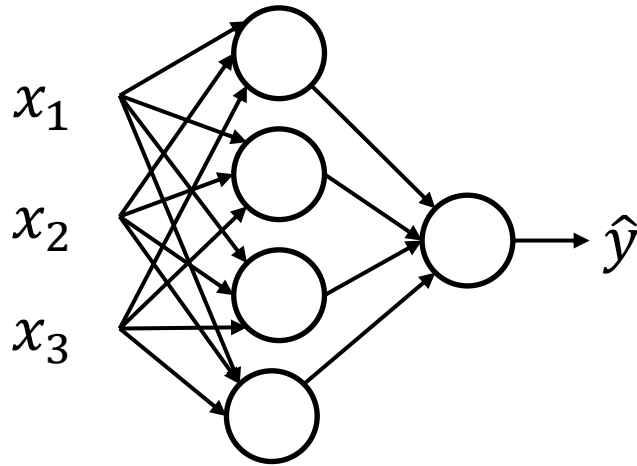
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

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

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

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

For loop across multiple examples

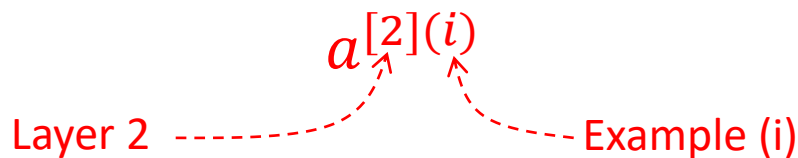


$$\begin{aligned}z^{[1]} &= W^{[1]}x + b^{[1]} \\a^{[1]} &= \sigma(z^{[1]}) \\z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} \\a^{[2]} &= \sigma(z^{[2]})\end{aligned}$$

$$\begin{array}{lcl}x & \text{----->} & a^{[2]} = \hat{y} \\x^{(1)} & \text{----->} & a^{[2](1)} = \hat{y}^{(1)} \\x^{(2)} & \text{----->} & a^{2} = \hat{y}^{(2)} \\... & \dots & ... \\x^{(m)} & \text{----->} & a^{[2](m)} = \hat{y}^{(m)}\end{array}$$

for i=1 to m:

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



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)

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

← # training examples →

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

↑ # hidden units ↓

Justification for vectorized implementation

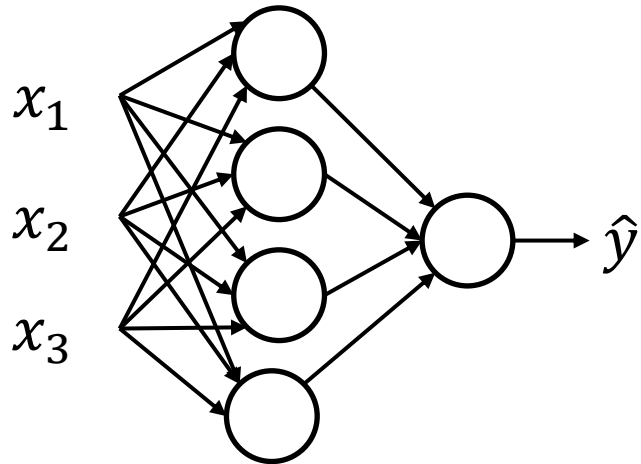
$$z^{1} = W^{[1]}x^{(1)} + b^{[1]} \quad , \quad z^{[1](2)} = W^{[1]}x^{(2)} + b^{[1]} \quad , \quad z^{[1](3)} = W^{[1]}x^{(3)} + b^{[1]}$$

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

X $Z^{[1]}$

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

Recap of vectorizing across multiple examples



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

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

for $i=1$ to m :

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

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

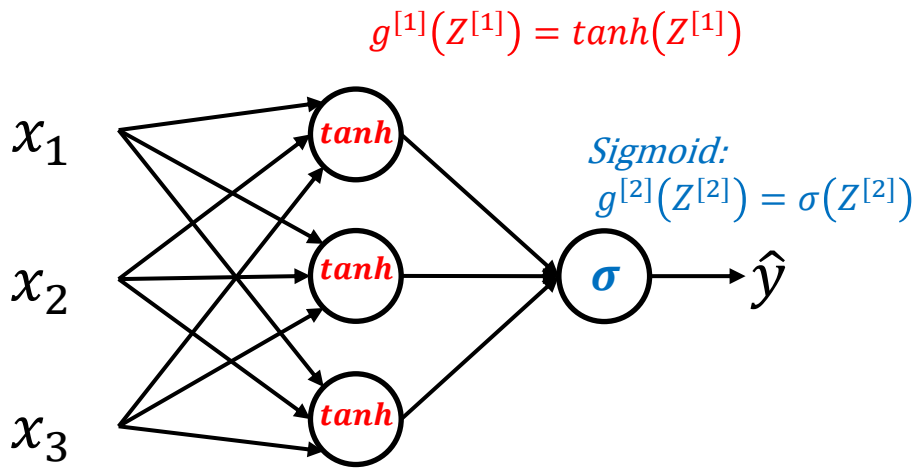
$$z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$$

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

$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} \\ A^{[1]} &= \sigma(Z^{[1]}) \\ Z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} &= \sigma(Z^{[2]}) \end{aligned}$$

A blue dashed arrow points from the circled X in the first equation to $A^{[0]}$.

Activation functions



Given X :

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

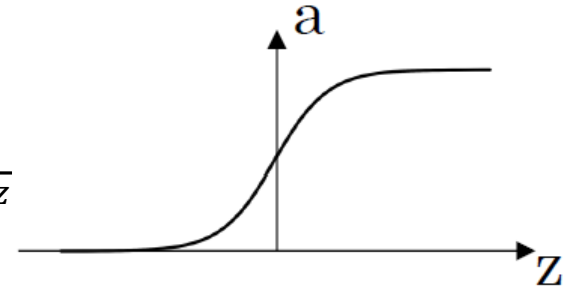
$$a^{[1]} = g^{[1]}(z^{[1]})$$

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

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

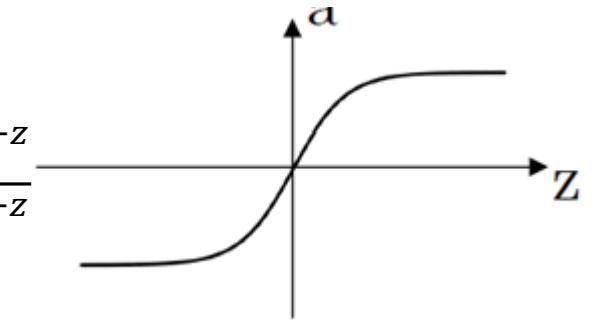
Sigmoid:

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



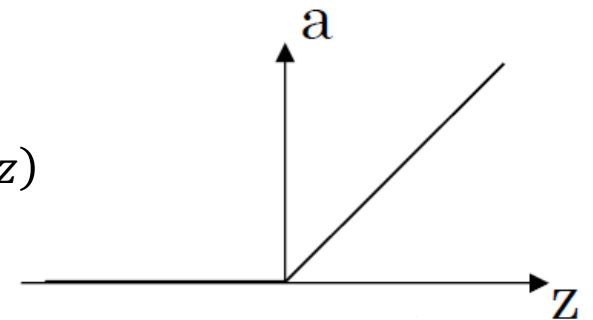
tanh:

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



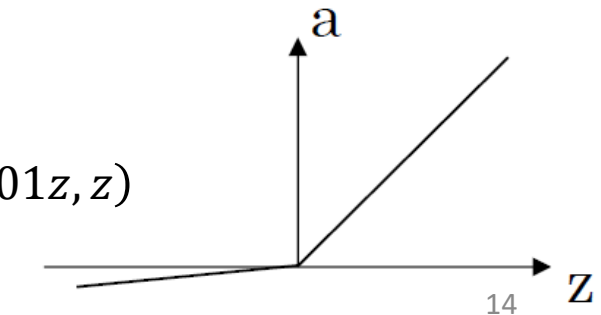
ReLU:

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



Leaky ReLU:

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

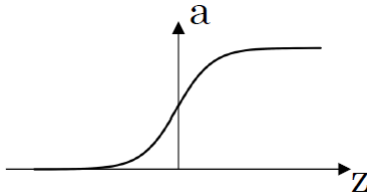


Pros and cons of activation functions

Sigmoid activation function:

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

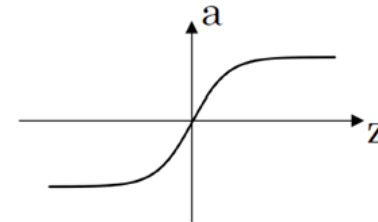
- Never use this, except for the output layer.
- if you are doing binary classification, or maybe almost never use this.



tanh activation function:

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

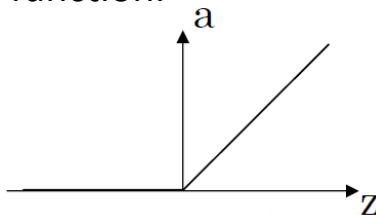
- The tanh is much strictly superior.



ReLU activation function:

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

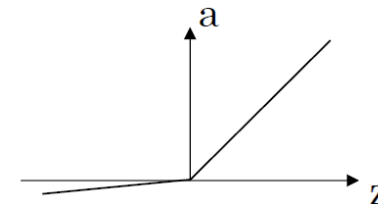
- The default and the most commonly used activation function is the ReLU.
- So if you're not sure what else to use, use the ReLU function.



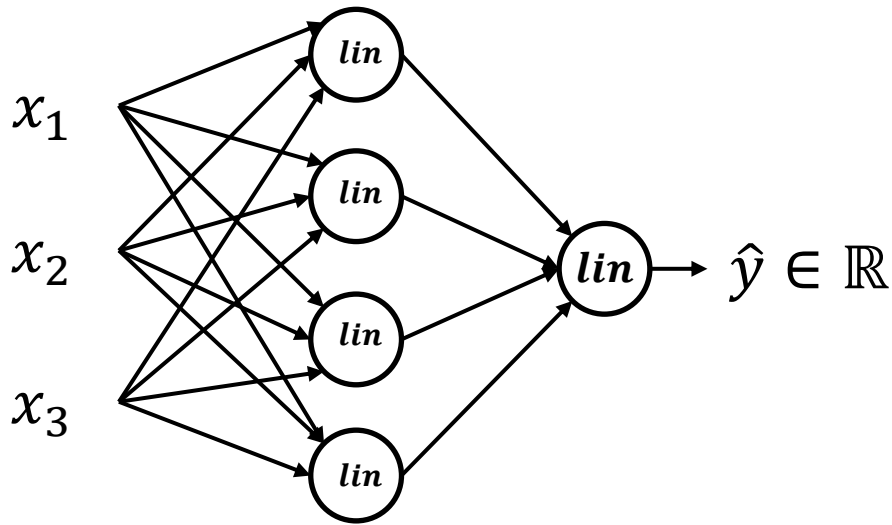
Leaky ReLU activation function:

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

- You can also try the leaky ReLU function.



Why do you need non linear activation functions?



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

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

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

Given X :

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

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

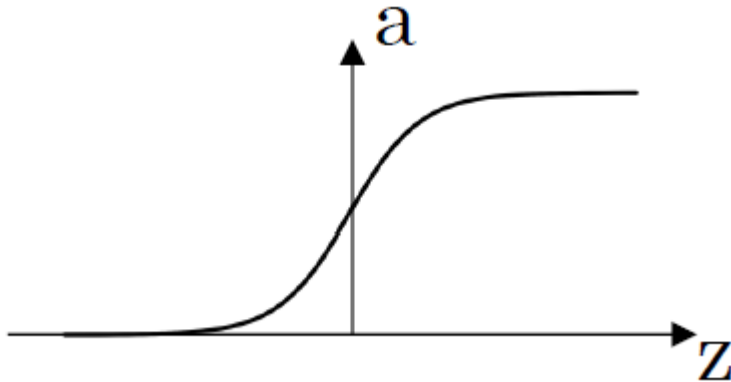
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

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

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

Derivatives of activation functions

Sigmoid activation function:



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

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

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

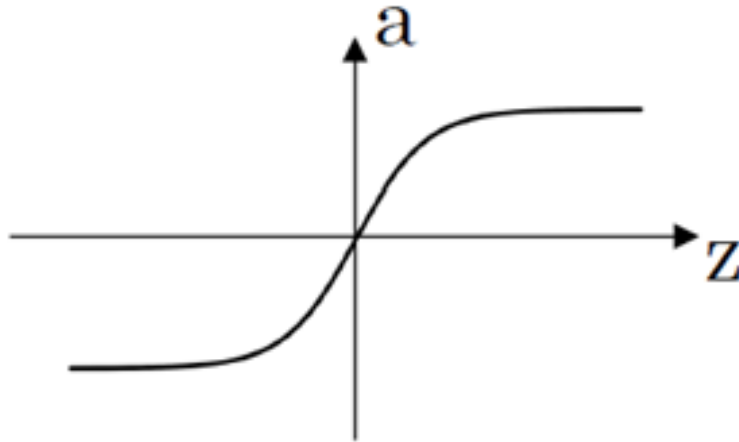
$$a = g(z), \quad g'(z) = a(1 - a)$$

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

- $z = 10 \Rightarrow g(z) \approx 1$
 $\Rightarrow \frac{d}{dz} g(z) \approx 1(1 - 1) \approx 0$
- $z = -10 \Rightarrow g(z) \approx 0$
 $\Rightarrow \frac{d}{dz} g(z) \approx 0(1 - 0) \approx 0$
- $z = 0 \Rightarrow g(z) = 1/2$
 $\Rightarrow \frac{d}{dz} g(z) = \frac{1}{2} \left(1 - \frac{1}{2} \right) = 1/4$

Derivatives of activation functions

Tanh activation function:



$$g'(z) = \frac{d}{dz}g(z) = \text{slope of } g(z) \text{ at } z$$
$$= 1 - (\tanh(z))^2$$

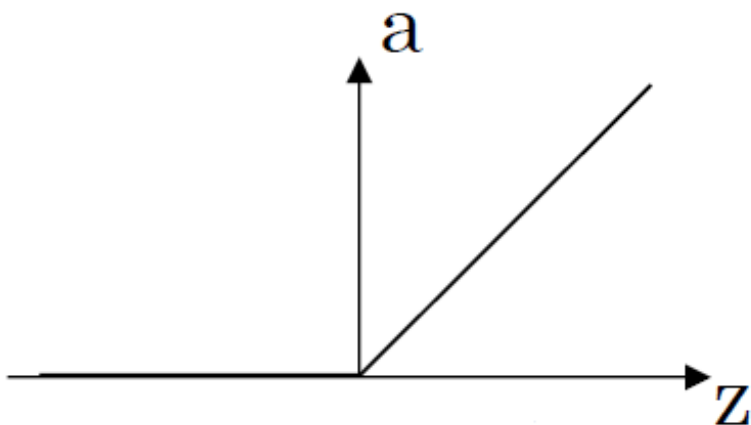
$$a = g(z), \quad g'(z) = 1 - a^2$$

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

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

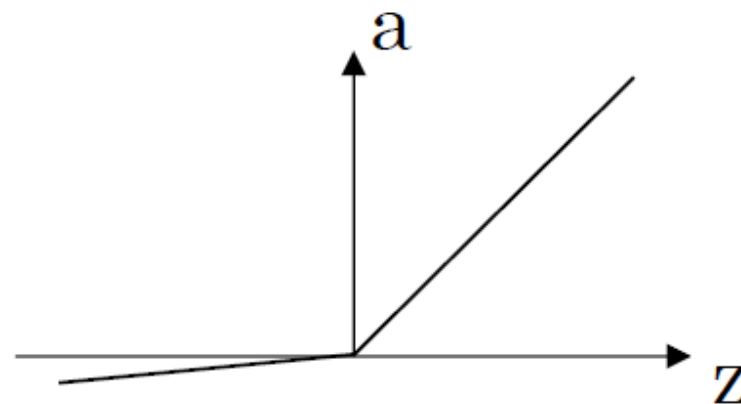
Derivatives of activation functions

ReLU and Leaky ReLU :



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

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



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

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

Gradient descent for neural networks

Parameters :

$$\begin{array}{cccc} W^{[1]} & b^{[1]} & W^{[2]} & b^{[2]} \\ (n^{[1]}, n^{[0]}) & (n^{[1]}, 1) & (n^{[2]}, n^{[1]}) & (n^{[2]}, 1) \end{array}$$

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

Cost function:

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

Gradient descent:

Repeat {

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

$$dW^{[1]} = \frac{\partial J}{\partial W^{[1]}}, db^{[1]} = \frac{\partial J}{\partial b^{[1]}}$$

$$W^{[1]} = W^{[1]} - \alpha dW^{[1]}$$

$$b^{[1]} = b^{[1]} - \alpha db^{[1]}$$

$$dW^{[2]} = \frac{\partial J}{\partial W^{[2]}}, db^{[2]} = \frac{\partial J}{\partial b^{[2]}}$$

$$W^{[2]} = W^{[2]} - \alpha dW^{[2]}$$

$$b^{[2]} = b^{[2]} - \alpha db^{[2]}$$

}

Formulas for computing derivatives

Forward propagation:

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

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

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

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

Back propagation:

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

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

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

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

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

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

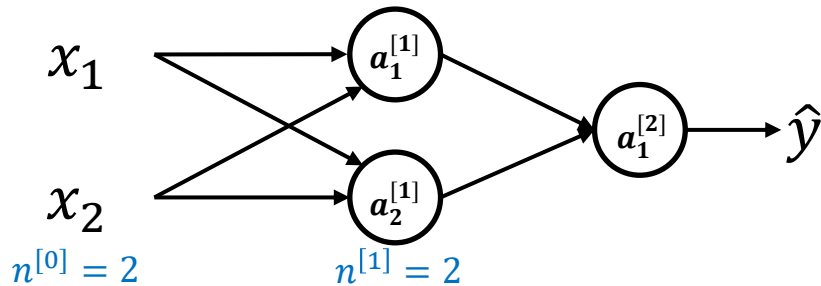
Element wise product

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

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

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

What happens if you initialize weights to zero?



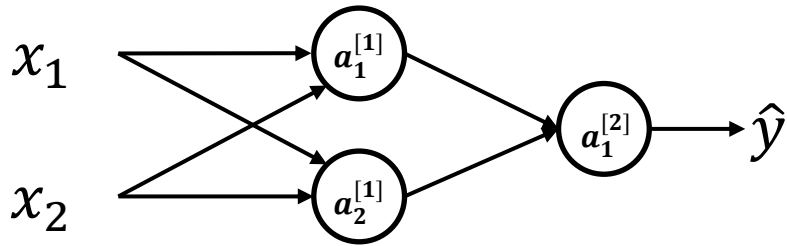
$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

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

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

- 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



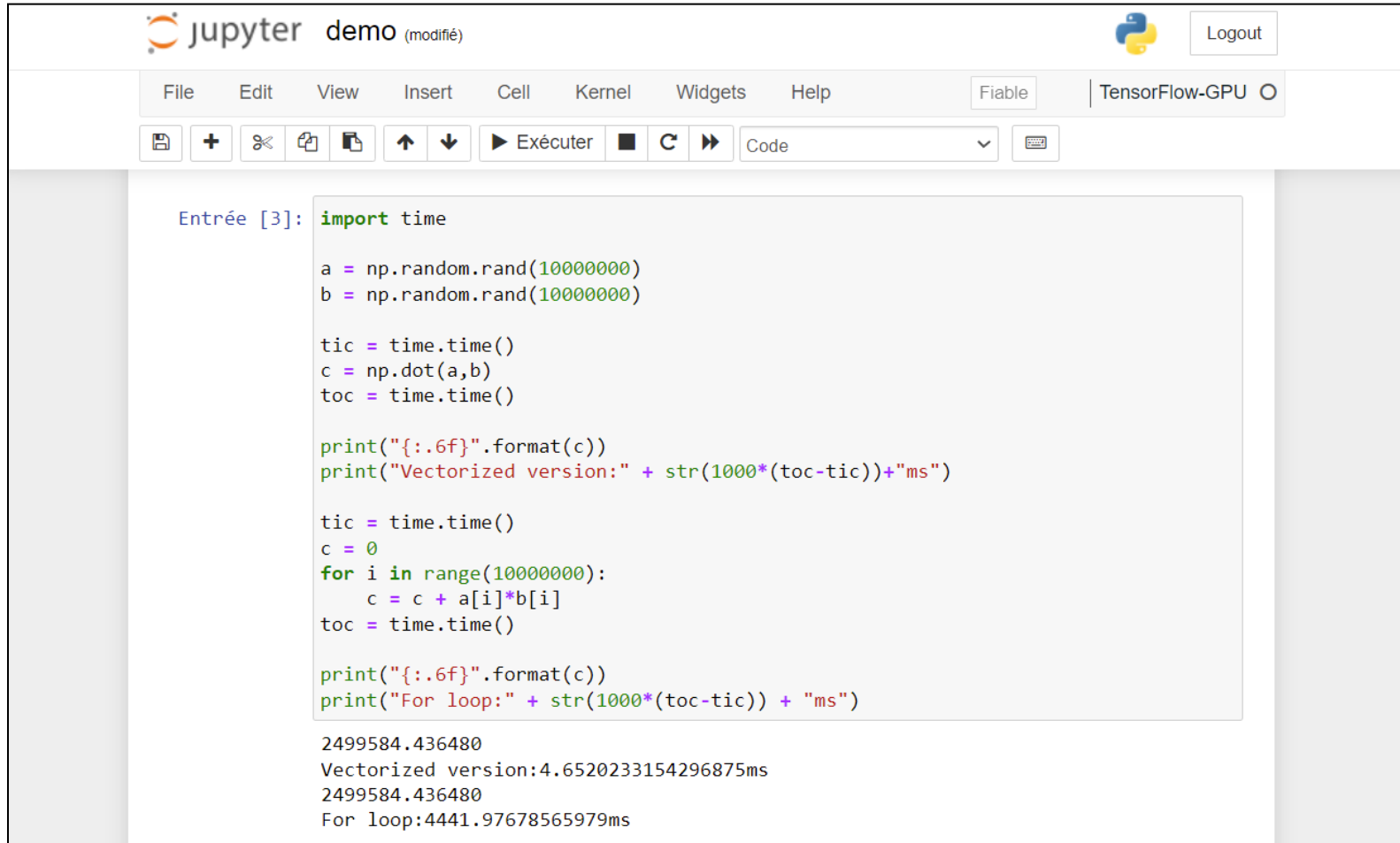
$$W^{[1]} = np.random.randn((2,2)) * 0.01$$

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

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

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

Vectorization demo



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

- Header: "jupyter demo (modifié)" on the left, Python logo and "Logout" button on the right.
- Menu bar: File, Edit, View, Insert, Cell, Kernel, Widgets, Help, Fiable, TensorFlow-GPU.
- Toolbar: Save, Add, Undo, Copy, Paste, Up, Down, Run (Exécuter), Stop, Refresh, Next, Code dropdown, and Help.
- Code cell content:

```
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")
```
- Output below the code cell:

```
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.