# MACHINE LEARNING

Universidad Francisco de Vitoria
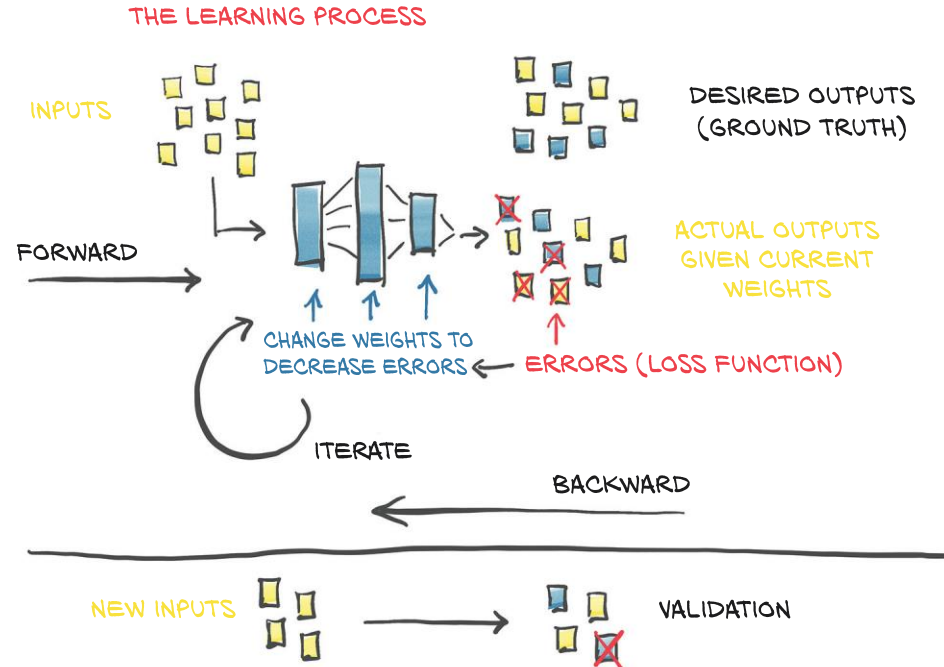
UFV Madrid

Forward pass

Backward pass

Optimization

Repetition Forward pass

**We will move from a Linear Regression to Neural Network**

Feature 1

Feature n

Output

Node 1

Feature 1

Node 2

Feature 2

Node n

Feature n

Output

We will go through this evolution using a binary classification problem

**Problem: to classify pictures as "picture" or "Non-Cat" picture**



1 (cat)

0 (cat)

y (label)

**How is a picture represented in a computer ?**



**RGB (Red, Green, Blue)**

64 pixels

64 pixels

**3 Arrays (64x64)**

x (feature)

Universidad
Francisco de
Vitoria

**UFV** Madrid

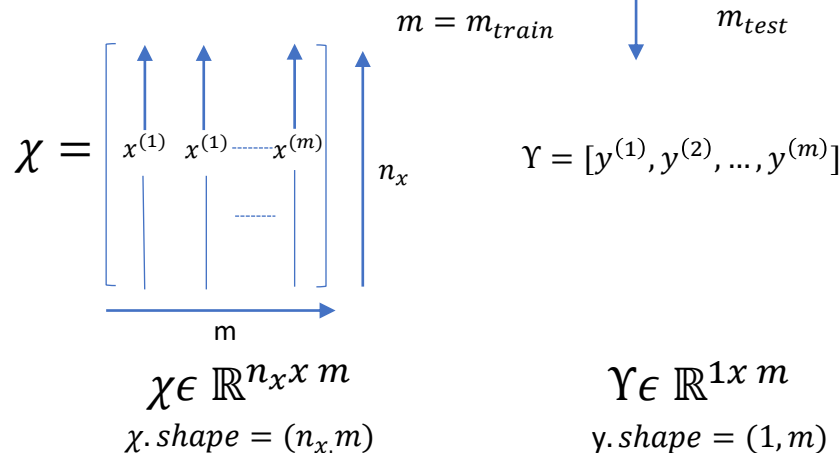## We will format our 3-dimensional array into a single vector x



x Vector with $64 \cdot 64 \cdot 3 = 12288\ components$

y Vector with just $1 component$

### Notation

Given a sample (x,y) $\longrightarrow$ $x \in \mathbb{R}^{n_x}, y \in \{0,1\}$

$$m\ training\ samples: \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

$$m = m_{train} \qquad m_{test}$$

$$\chi = \begin{bmatrix} x^{(1)} & x^{(1)} & \cdots & x^{(m)} \\ & & \cdots & \end{bmatrix} \Big\} n_x$$

$$\Upsilon = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$

m

$$\chi \in \mathbb{R}^{n_x\ x\ m}$$
$$\chi.\,shape = (n_x, m)$$

$$\Upsilon \in \mathbb{R}^{1\ x\ m}$$
$$y.\,shape = (1, m)$$

Universidad
Francisco de
Vitoria

**UFV** Madrid

## A logistic regression is a Neural Network used for binary classification

**Our problem:** Given a sample x, we want to $\hat{y}$ as the probability of being a cat $\quad \hat{y} = P(y = 1|x) \quad 0 \leq \hat{y} \leq 1$

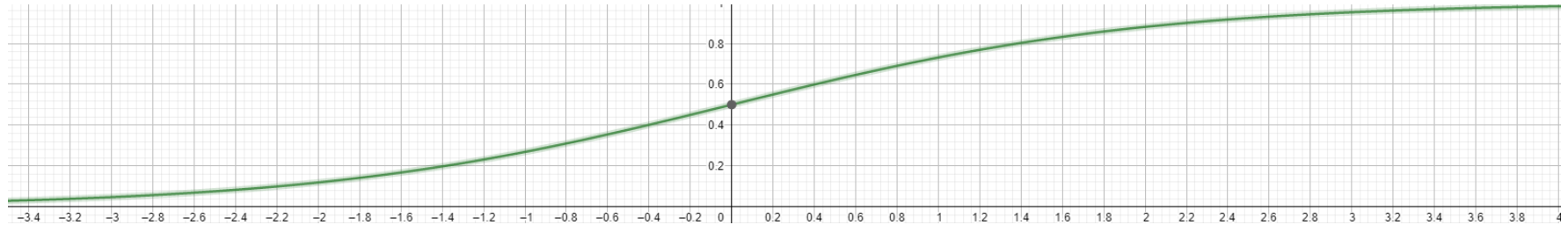**Our model:** $z = w^T \cdot x + b, where \; w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$



**But... , this function returns values greater than 1 !!!**

**And.., probability should be a real number between, and 1**

Universidad Francisco de Vitoria
UFV Madrid



$We\ call\ z = w^T + b$

$$\hat{y} = \ \sigma(z) \quad \sigma(z) = \frac{1}{1 + e^{-z}}$$

If z is large + …

$$\sigma(z) \approx \frac{1}{1 + 0} \approx 1$$

Now our output is number between 0 and 1

If z is large -…

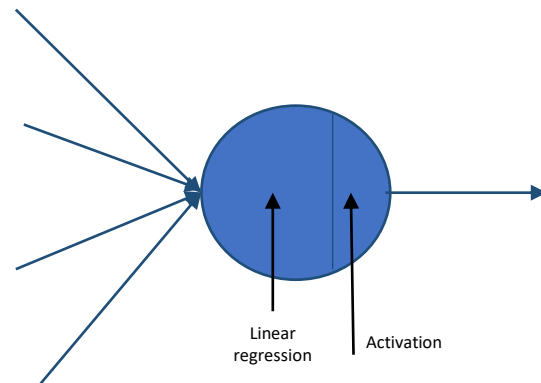$$\sigma(z) \approx \frac{1}{1 + big} \approx 0$$

A neuron in Machine Learning is a function composition between:

- linear equation (regression function)

- and a non-linear equation (activation function)

$$\hat{y} = g \circ f(w, b) = g(f(w, b))$$



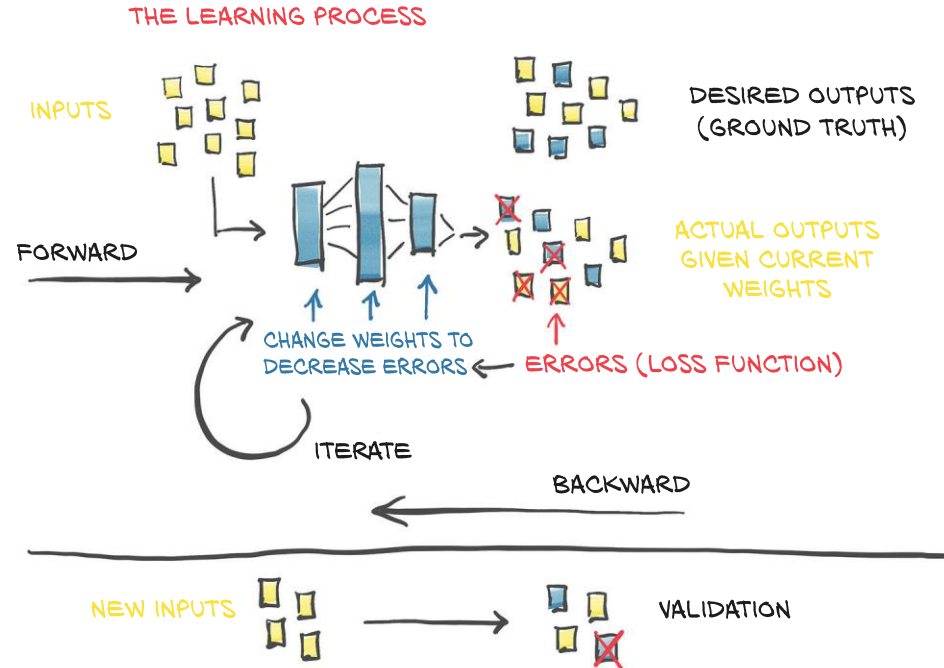Linear
regression

Activation

**In our Cat classifier (binary classification):**

$$\hat{y} = \sigma \circ z(w, b) = \sigma\big(z(w, b)\big) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-(w^T x + b)}}$$

$w^T \ is \ vector, b \ is \ a \ real \ number$

✓ Forward pass

➡ Backward pass

Optimization

Repetition Forward pass

THE LEARNING PROCESS

INPUTS

DESIRED OUTPUTS
(GROUND TRUTH)

FORWARD

ACTUAL OUTPUTS
GIVEN CURRENT
WEIGHTS

CHANGE WEIGHTS TO
DECREASE ERRORS ← ERRORS (LOSS FUNCTION)

ITERATE

BACKWARD

NEW INPUTS → VALIDATION
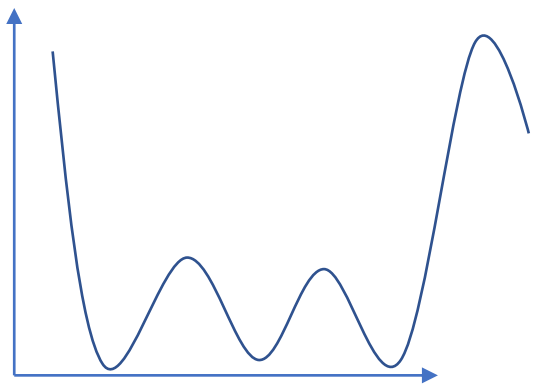
**Logistic Regression Cost Function**

i is the i-th example

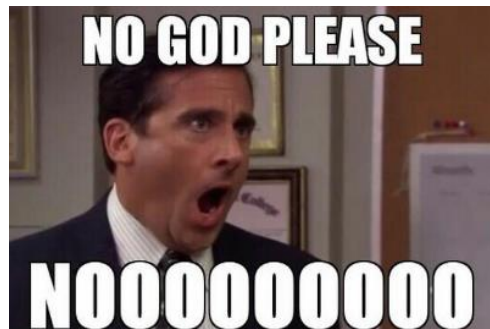$$Given\ \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\},$$

we want $\hat{y}^{(i)} \approx y^{(i)}$

$$Our\ cost\ function\ should\ be\ \mathcal{L}(\hat{y}, y) = (\hat{y} - y)^2$$

Local minimum issue

NO GOD PLEASE

NOOOOOOOOO

Gradient Descent doesn't work well with local minimum

**Specific Loss Function for Logistic Regression**

Loss function $\longrightarrow$ $\mathcal{L} = (\hat{y}, y) - (y \cdot \log(\hat{y}) - (1 - y) \cdot \log(1 - \hat{y}))$

$If\ y = 1, \mathcal{L}(\hat{y}, y) = -\log(\hat{y})$ $\qquad If\ \hat{y}\ close\ to\ 1, \mathcal{L} \approx 0$

$If\ y = 0, \mathcal{L}(\hat{y}, y) = -\log(1 - \hat{y})$ $\qquad If\ \hat{y}\ close\ to\ 0, \mathcal{L} \approx 0$

$$\Im(w, b) = \frac{1}{m}\sum_{i=1}^{m}\mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m}\sum_{i=1}^{m}[y^{(i)}\log\hat{y}^{(i)} + (1 - y^{(i)})\log(1 - \hat{y}^{(i)})]$$

Cost function

**Loss Function versus Cost Function**

Loss function measures de accuracy in a single example

$$\mathcal{L}(\hat{y}, y) = -(y \cdot \log(\hat{y}) - (1 - y) \cdot \log(1 - \hat{y}))$$

Cost function measures de accuracy on the whole dataset

$$\mathfrak{J}(w, b) = -\frac{1}{m} \sum_{i=1}^{m} [y^{(i)} log\hat{y}^{(i)} + (1 - y^{(i)})\log(1 - \hat{y}^{(i)})]$$

**We're going to suppose that we only have 1 sample with two features**

$x_1$
$w_1$

$x_2$
$w_2$

$b$

$$z = w_1 \cdot x_1 + w_2 \cdot x_2 + b \quad \Rightarrow \quad a = \sigma(z) \quad \Rightarrow \quad \mathcal{L}(a, y)$$

First backward step

Second backward step

$$\frac{d\mathcal{L}(a, y)}{da}$$

Third backward step

$$\frac{d\mathcal{L}(a, y)}{dz}$$

Ground truth

$$[\frac{\partial \mathcal{L}(a,y)}{\partial w_1}, \frac{\partial \mathcal{L}(a,y)}{\partial w_w}, \frac{\partial \mathcal{L}(a,y)}{\partial b}]$$
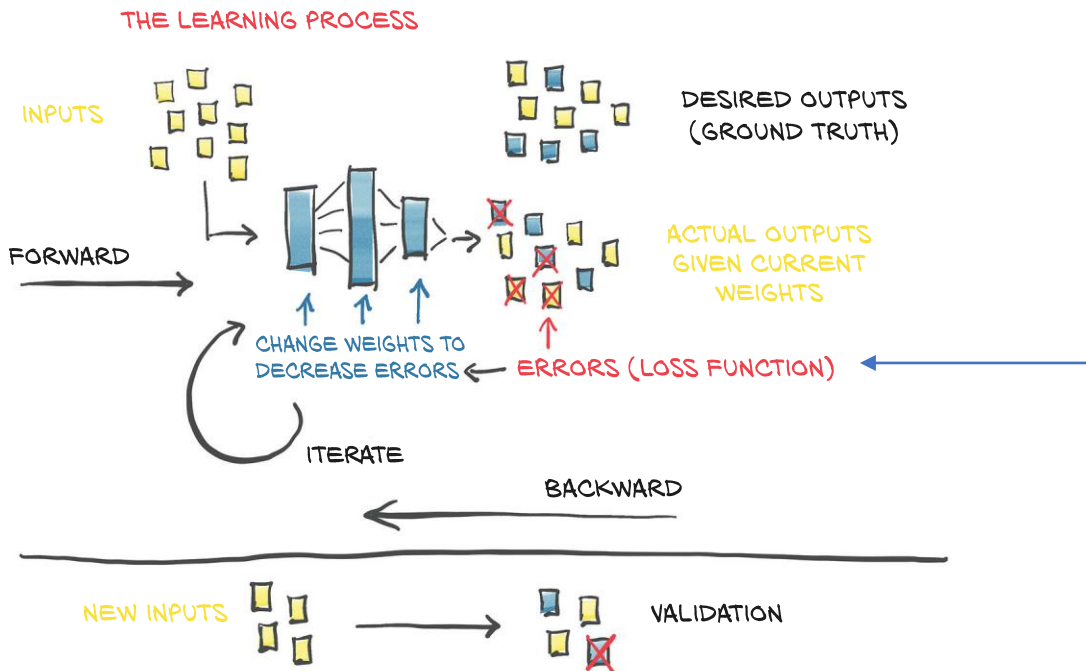
output

feature

$$\frac{\partial \mathcal{L}(a,y)}{\partial w_1} = \frac{dz}{dw_1} \cdot \frac{da}{dz} \cdot \frac{d\mathcal{L}(a,y)}{dy} = x_1 \cdot a \cdot (1 - a) \cdot \left( -\frac{y}{a} + \frac{1 - y}{1 - a} \right) = x_1(a - y)$$

**Therefore, the gradinet vector for the cost function is**

$$\nabla \mathcal{L}[w_1, w_2, b] = [x_1(a - y), x_2(a - y), (a - y)]$$

✓ Forward pass

✓ Backward pass

➡ Optimization

Repetition Forward pass



THE LEARNING PROCESS

INPUTS

DESIRED OUTPUTS (GROUND TRUTH)

FORWARD

ACTUAL OUTPUTS GIVEN CURRENT WEIGHTS

CHANGE WEIGHTS TO DECREASE ERRORS ← ERRORS (LOSS FUNCTION)

ITERATE

BACKWARD

NEW INPUTS

VALIDATION

**Remind the gradient descent algorithm used in the linear regression practice**
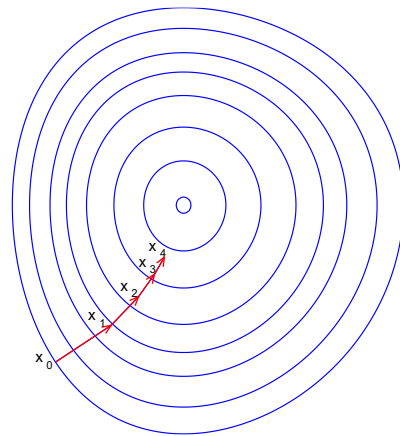
$$z_{n+1} = z_n - \gamma \nabla \mathcal{L}(z_n)$$

$$w_{1_{n+1}} = w_{1_n} - \gamma \nabla \mathcal{L}\left(w_{1_n}\right) = w_{1_n} - \gamma \cdot x_n(a_n - y_n)$$
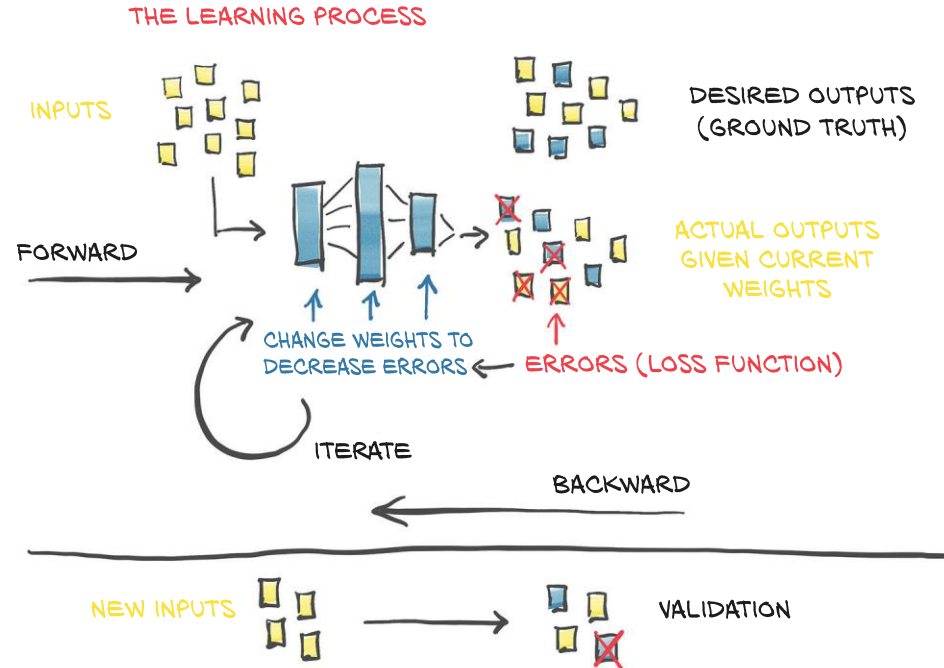
$$w_{2_{n+1}} = w_{2_n} - \gamma \nabla \mathcal{L}\left(w_{2_n}\right) = w_{2_n} - \gamma \cdot x_n(a_n - y_n)$$

$$b_{n+1} = b_n - \gamma \nabla \mathcal{L}(b_n) = b_n - \gamma \cdot (a_n - y_n)$$

$\gamma$ = **Learning rate**

- Forward pass

- Backward pass

- Optimization

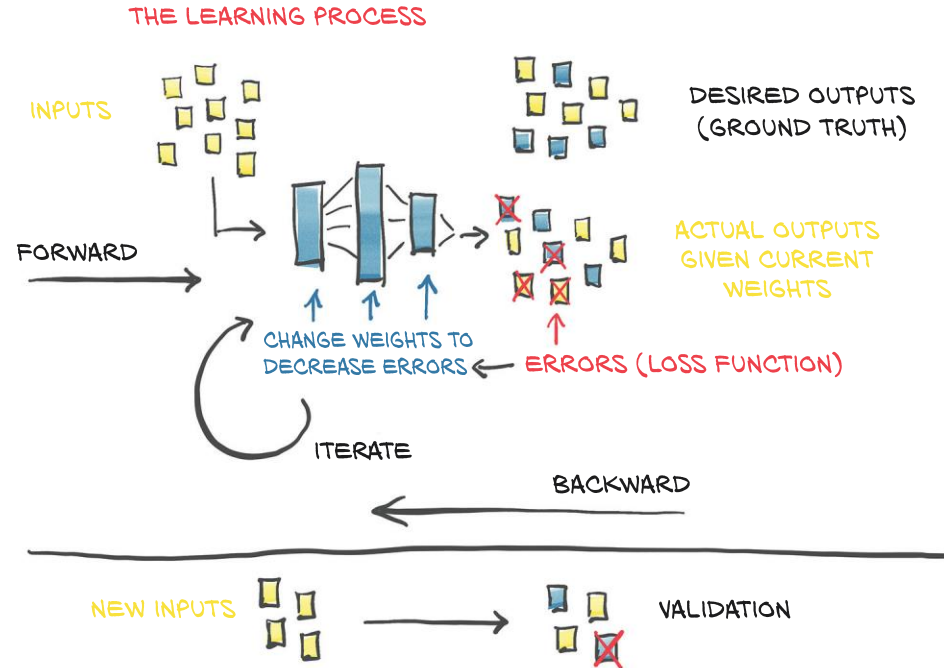- Repetition Forward pass

**We will use a for loop to iterate as many times as we consider (epochs)**

```
for epochs in range(1,n_epochs+1)

    a = compute our output
    Loss = compute our loss
    Gradient = compute our gradient
    Params in epoch n = (params in epoch n-1)- learning rate x Gradient Descent
```

- ✓ Forward pass

- ✓ Backward pass

- ✓ Optimization

- ✓ Repetition Forward pass

How will our algorithm be with m examples instead of just

1 example?

Universidad
Francisco de
Vitoria

**UFV** Madrid

**What happens if I have to compute the algorithm for m examples instead of 1 example??**

**Let's review our trainin loop**

```
for epoch in epochs:
    for i in range(1,m+1):
```

$$z^{(i)} = w^T \cdot x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$   ⟵ **Forward Function**

$$\mathcal{L} = \mathcal{L} + (-[y^{(i)} - \log a^{(i)} + (1 - y^{(i)})\log(1 - a^{(i)})])$$   ⟵ **Loss Function**

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$db = db + dz^{(i)}$$   ⟵ **Gradient Vector**

```
        for n in range(1,n_features):
```

$$dw_n = x_n^{(i)} \cdot dz^{(i)}$$

$$\mathfrak{I} = \frac{\mathcal{L}}{m}, \frac{dw_1}{m}, \dots, \frac{dw_n}{m}, \frac{db}{m}$$   ⟵ **Cost Function**

```
        for n in range(1,n_features):
```
⟵ **Optimization**

$$w_n = w_n - \gamma \cdot dw_n$$

$$o(n^3)$$

# MACHINE LEARNING

## VECTORIZATION

$o(n^3)$ **algorithms don't escalate, moreover in Machine Learning environments, where data tends to grow exponentially**

**Non-Vectorized**

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

$z = 0$

for i in range (1,n_features+1)

   $z = z + w[i] \cdot x[i]$

**Vectorized**

$z = np.dot(w, x)$

$w^T x$

```
# -*- coding: utf-8 -*-
"""
Created on Thu Oct 20 10:59:15 2022

@author: PedroAnquela
"""

import numpy as np
import time

dim = 10000000

w = np.random.rand(dim)
x = np.random.rand(dim)

start = time.time()

a = np.dot(w,x)

end = time.time()

print("Vectorized version: {0:.4f} ms.".format(1000*(end-start)))

######################################
# Non-Vectorized Version             #
######################################
c = 0
start = time.time()

for i in range (dim):
    c += w[i] * x[i]

end = time.time()

print("Non-Vectorized version: {0:.4f} ms.".format(1000*(end-start)))
```

```
In [8]: runfile('C:/Users/PedroAnquela/Documents/Projectos_Local/UFV/
prog_ii_2020/ml/vectorization/vectorization.py', wdir='C:/Users/
PedroAnquela/Documents/Projectos_Local/UFV/prog_ii_2020/ml/
vectorization')
Vectorized version: 9.9742 ms.
Non-Vectorized version: 3460.2816 ms.
```

**384 times
faster
when vectorizing**

**We turn data into numpy vectors**

$$A = [a^{(1)}, .. a^{(m)}]$$

$$Y = [y^{(1)}, .. y^{(m)}]$$

$$dZ = A - Y$$

$$db = \frac{1}{m} \sum_{i=1}^{m} dz^{(i)} = \frac{1}{m} np.sum(dZ)$$

$$dw = \frac{1}{m} X \cdot dz^T = \frac{1}{m} [x^{(1)} \dots x^{(m)}] \begin{bmatrix} dz^{(1)} \\ \dots \\ dz^{(m)} \end{bmatrix}$$

$$(1, m) \cdot (m, 1)$$

**Let's review our trainin loop**

$J = 0, dw_1 = 0, dw_2 = 0, db = 0$

```
for epoch in epochs:
    for i in range(1,m+1):
```
$$z^{(i)} = w^T \cdot x^{(i)} + b$$
$$a^{(i)} = \sigma(z^{(i)})$$
$$\mathcal{L} = \mathcal{L} + (-[y^{(i)} - \log a^{(i)} + (1 - y^{(i)})\log(1 - a^{(i)})])$$
$$dz^{(i)} = a^{(i)} - y^{(i)}$$
$$db = db + dz^{(i)}$$
```
        for n in range(1,n_features):
```
$$dw_n = x_n^{(i)} \cdot dz^{(i)}$$
$$\mathfrak{I} = \mathcal{L}/m, dw_1/m, \dots, dw_n/m, db/m$$
```
        for n in range(1,n_features):
```
$$w_n = w_n - \gamma \cdot dw_n$$

```
for epoch in epochs:
```

$$Z = w^T X + b = np.dot(w.T, X) + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

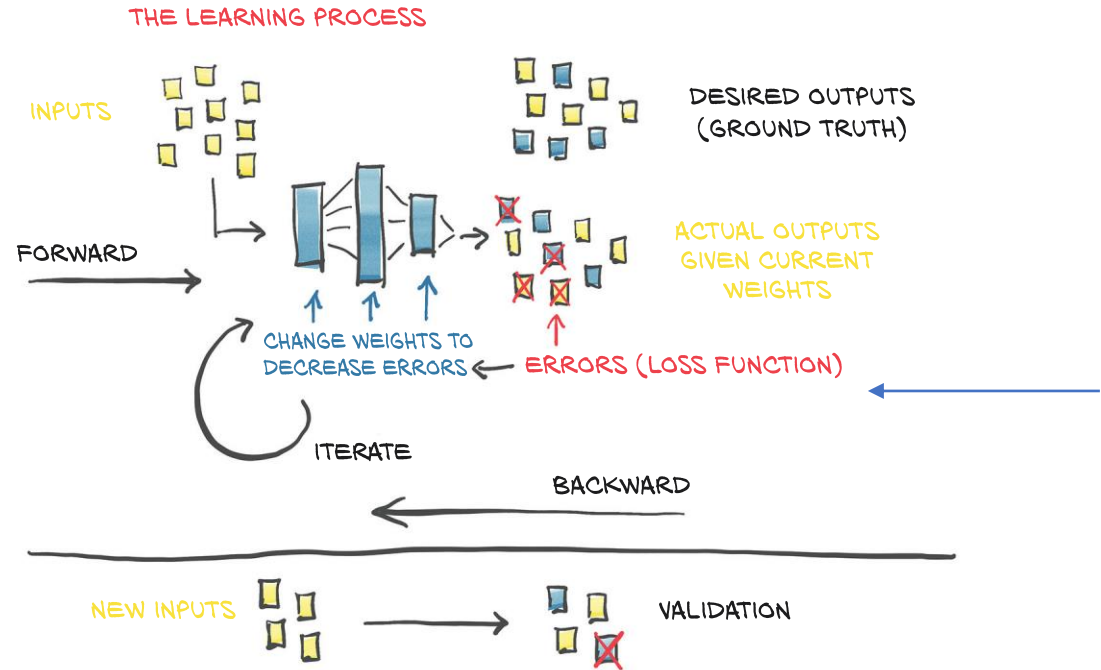$$db = \frac{1}{m} \sum_{i=1}^{m} \frac{1}{m} np.sum(dZ)$$

$$dw = \frac{1}{m} X \cdot dz^T$$

$$w = w - \gamma dw$$

$$b = b - \gamma db$$

✓ Forward pass

✓ Backward pass

✓ Optimization

✓ Repetition Forward pass

THE LEARNING PROCESS

INPUTS

DESIRED OUTPUTS
(GROUND TRUTH)

FORWARD

ACTUAL OUTPUTS
GIVEN CURRENT
WEIGHTS

CHANGE WEIGHTS TO
DECREASE ERRORS ← ERRORS (LOSS FUNCTION)

ITERATE

BACKWARD

NEW INPUTS

VALIDATION

**Remembering a basic concept of Linear Algebra**

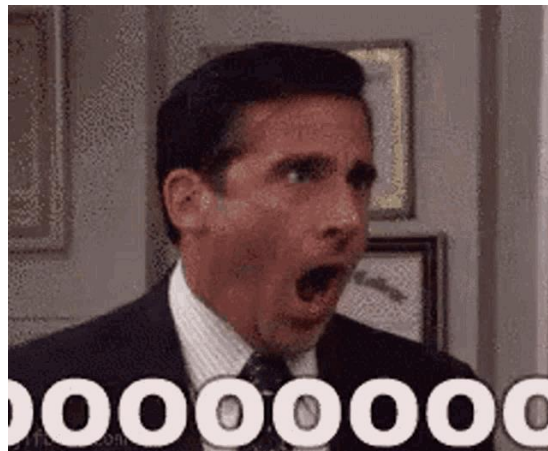$$\begin{bmatrix} 100 \\ 101 \\ 102 \\ 103 \end{bmatrix} + 100 =$$

Can I do this operation ?

Two matrices must have an equal number of rows and columns to be added.

In which case, the sum of two matrices **A** and **B** will be a matrix which has the same number of rows and columns as **A** and **B**.

The sum of **A** and **B**, denoted **A + B**, is computed by adding corresponding elements of **A** and **B**.

- Open your text editor
- Generate a numpy array with the following dimensions (4,1)
- add the number 100 to your numpy.array.
- Print the result
- Did it work? What was the result?

```python
import numpy as np

a = np.array([[100],
              [101],
              [102],
              [103]
              ])

print(a + 100)
```

**Numpy broadcasts non-adjusted shapes to allow the addition**

$$\begin{bmatrix} 100 \\ 101 \\ 102 \\ 103 \end{bmatrix} + 100 = numpy\ broadcasting = \begin{bmatrix} 100 \\ 101 \\ 102 \\ 103 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 200 \\ 201 \\ 202 \\ 203 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 5 \\ 3 & 4 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix} = numpy\ broadcasting = \begin{bmatrix} 1 & 2 & 5 \\ 3 & 4 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 305 \\ 103 & 204 & 306 \end{bmatrix}$$

**Numpy broadcasts non-adjusted shapes to allow the addition**

**Numpy broadcasts general principles**

$$A_{(m,n)} \; + \; B_{(1,n)} \longrightarrow A_{(m,n)} + B_{(m,n)}$$

$$A_{(m,n)} \; - \; C_{(m,1)} \longrightarrow A_{(m,n)} + C_{(m,n)}$$

$$A_{(m,n)} + n \; \in \; \mathbb{R} \longrightarrow A_{(m,n)} + B_{(m,n)} \,, b_{ij} = n \; \forall b_{ij} \in B$$

**Never let Numpy broadcast work in Matrix Multiplication**

$$A_{(m,n)} \cdot B_{(n,m)}$$

**Number of columns of first matrix HAS TO match the number of rows of the second matrix**

- Given the following matrix

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

- And the following vector

$$B = [2,3,4]$$

- Use numpy to calculate A*B

Use the **command** numpy.reshape to fix the dimensions you need.

## Never let numpy decide by itself dimensions of your matrix

```python
import numpy as np

a = np.ones(5)

print(a)
print(a.shape)


a = a.reshape(5,1)

print(a)
print(a.shape)

a = a.reshape(1,5)

print(a)
print(a.shape)
```

```
[1. 1. 1. 1. 1.]
(5,)
```

```
[[1.]
 [1.]
 [1.]
 [1.]
 [1.]]
(5, 1)
```

```
[[1. 1. 1. 1. 1.]]
(1, 5)
```

## Cat / non-Cat classifier

- Build a Logistic Regression (Shallow Neural Network) to recognize cats

- Classify your own photo as Cat or Non-Cat

- Loops are not permitted, unless there is an explicit function to do so

- Use Jupyter Notebook to fill up the Assignment (conda install jupyter)

- Clone the Notebook from the following link

- At the top of a Given Block you will find a #Graded Función tag, you will write you own code between the ###Start Solution Here ### and the ###End Solution Here### comments

- Take your time, thing about all the concepts explained in the previous lectures. Check out if your model has converged and it can recognize a cat