

Ph.D Juan David Martínez Vargas
Ph. D Raul Andrés Castañeda

Escuela de Ciencias Aplicadas e Ingeniería

Lecture 02

Deep Learning

Agenda

- Modelos neuronales
- Perceptron: Single layer NNs
- Álgebra lineal y cálculo para Deep Learning

Redes Neuronales: Inspiradas por el Cerebro

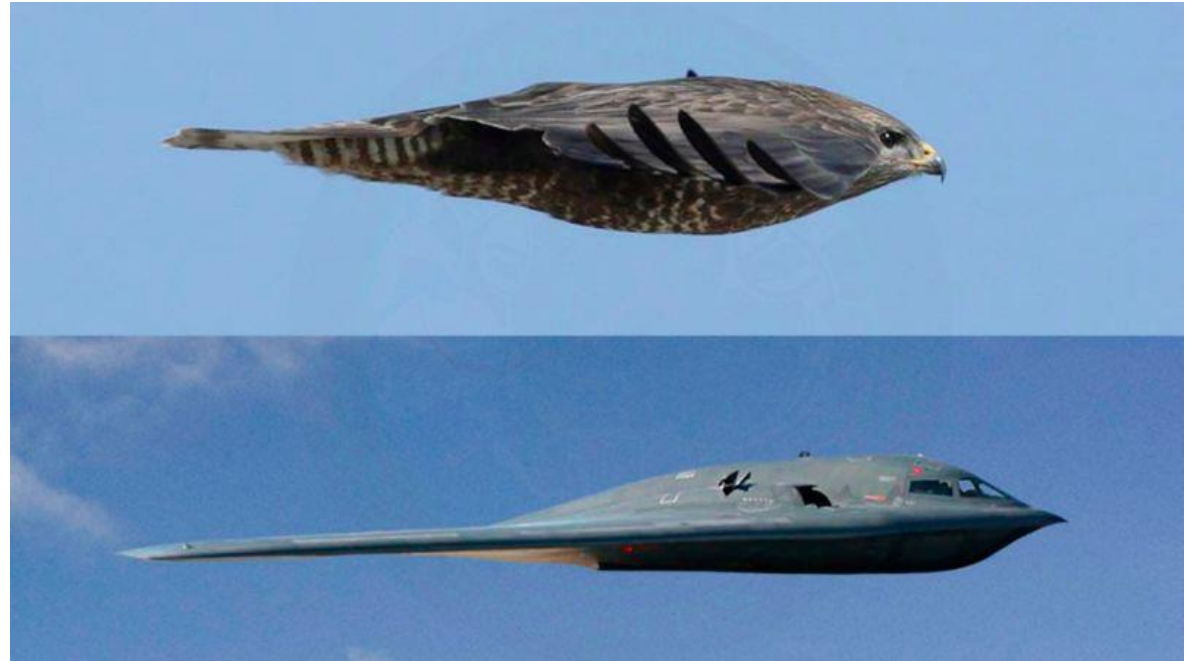


<https://www.verywellmind.com/how-brain-cells-communicate-with-each-other-2584397>

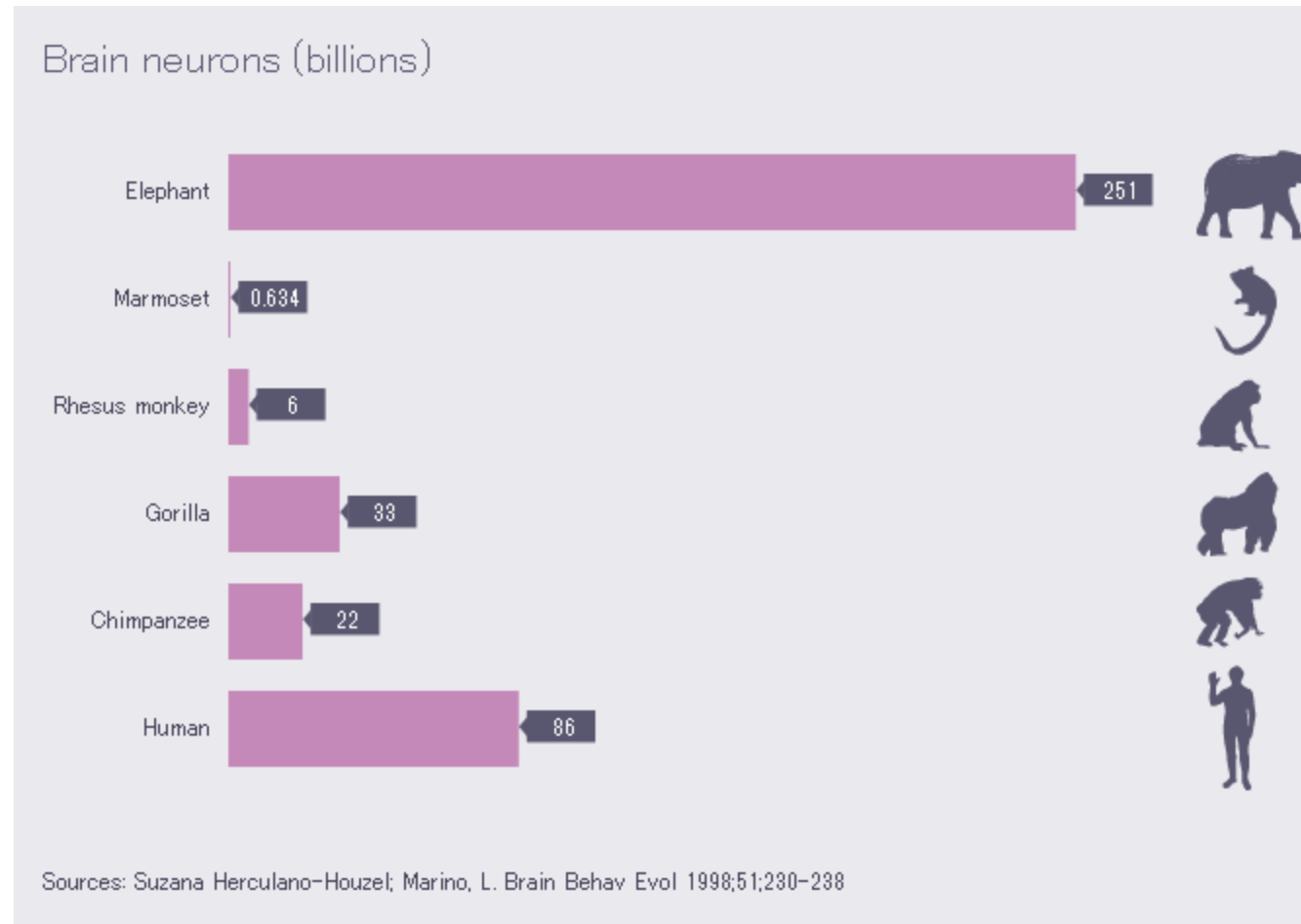
Redes Neuronales: Inspiradas por el Cerebro



<https://medium.com/@adsactly/is-it-a-bird-is-it-a-plane-biomimicry-in-airplanes-9862d331df2e>

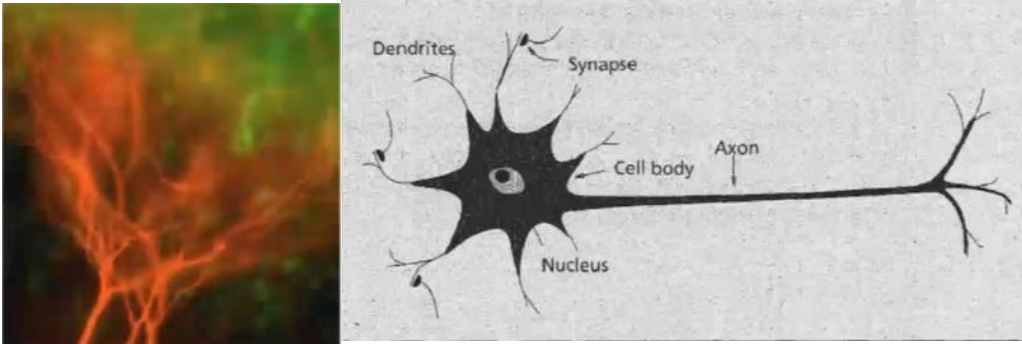


Número de neuronas en el cerebro



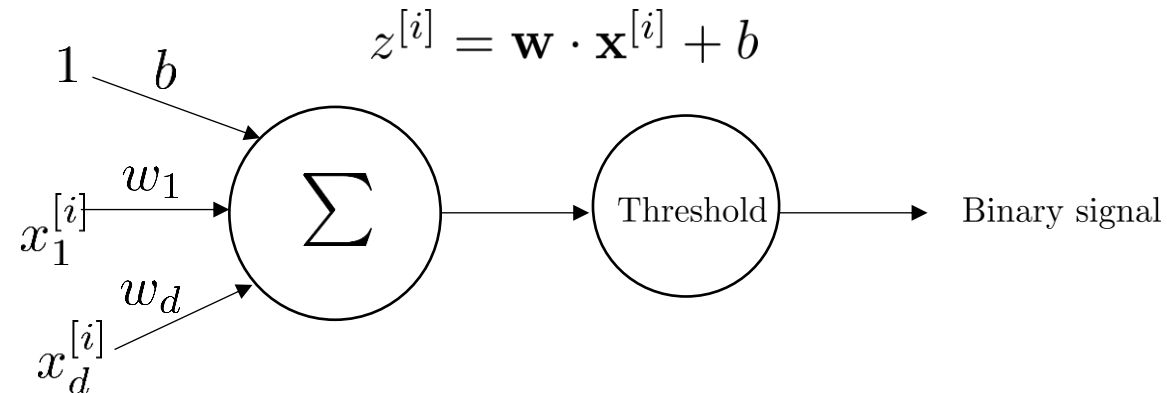
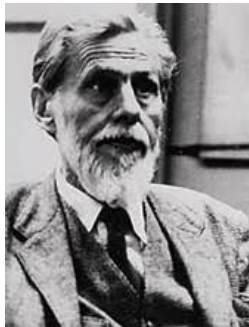
[https://en.wikipedia.org/wiki/List_of_animals_by_number_of_neurons#/media/File:Brain_size_comparison_-_Brain_neurons_\(billions\).png](https://en.wikipedia.org/wiki/List_of_animals_by_number_of_neurons#/media/File:Brain_size_comparison_-_Brain_neurons_(billions).png)

Modelo de Neurona de McCulloch Pitts



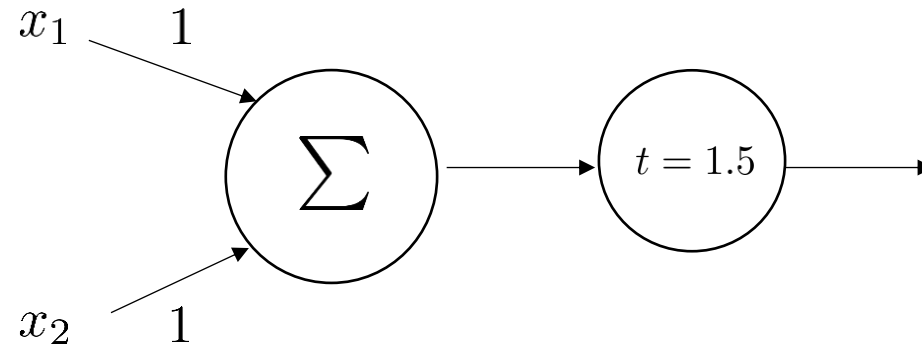
A LOGICAL CALCULUS OF THE IDEAS IMMANENT IN NERVOUS ACTIVITY

WARREN S. MCCULLOCH and WALTER H. PITTS 1943



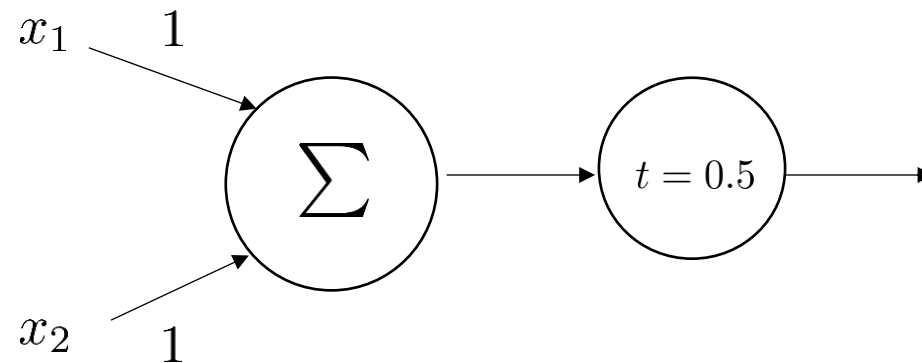
And

x_1	x_2	Out
0	0	0
0	1	0
1	0	0
1	1	1



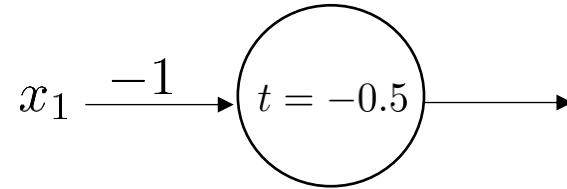
Or

x_1	x_2	Out
0	0	0
0	1	1
1	0	1
1	1	1



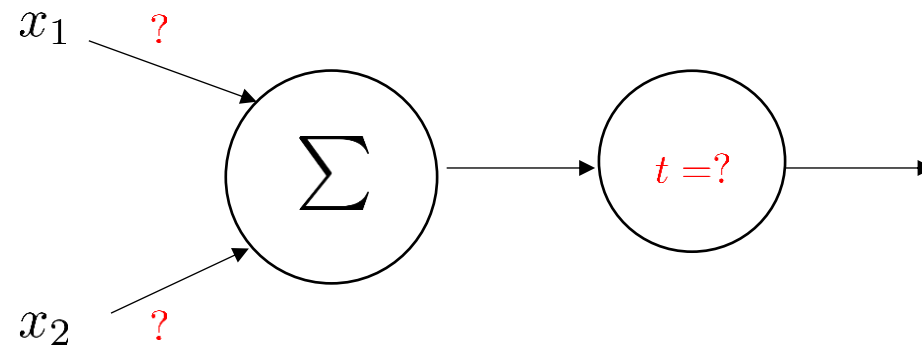
Not

x_1	Out
0	1
1	0

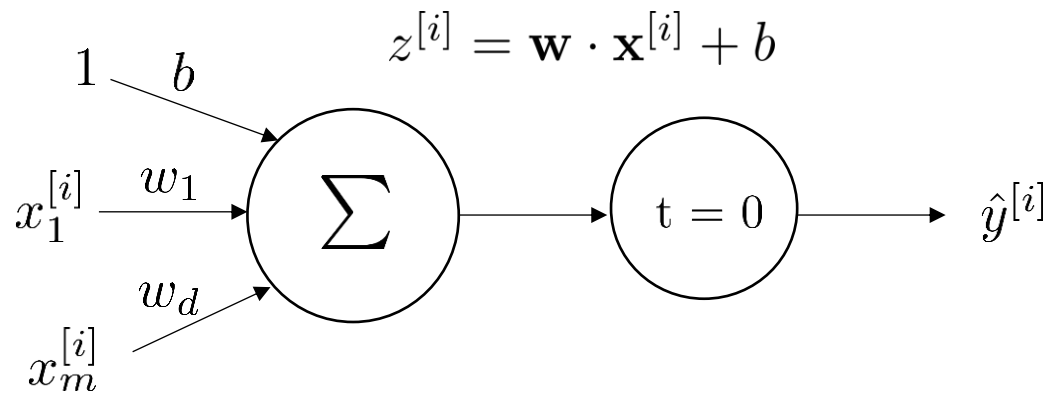


Xor

x_1	x_2	Out
0	0	0
0	1	1
1	0	1
1	1	0

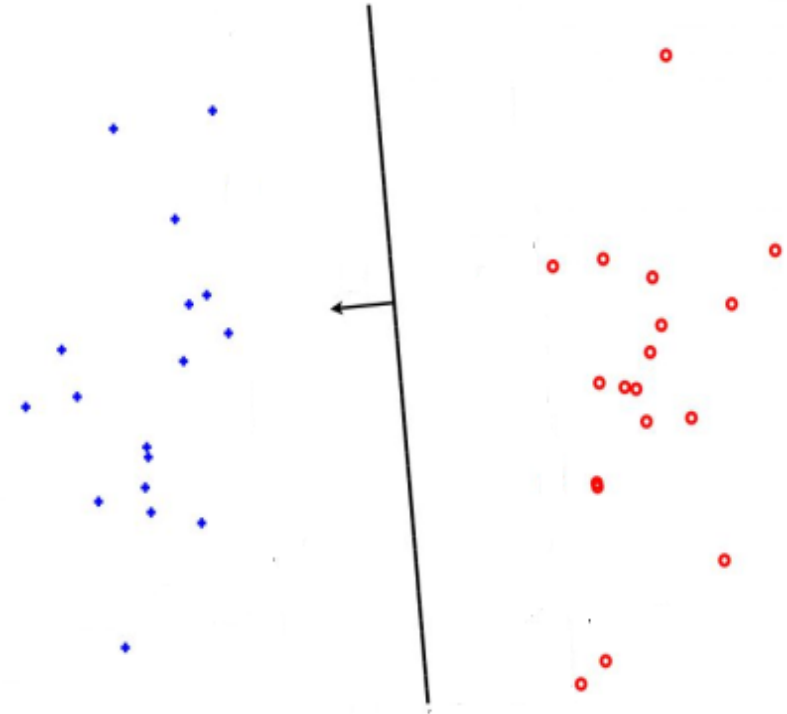


Perceptrón



$$\mathbf{x}^{[i]} = [1, x_1^{[i]}, \dots, x_m^{[i]}] \quad w_0 = b$$

$$z^{(i)} = \mathbf{w} \cdot \mathbf{x}^{[i]}$$



$$\mathcal{D} = \{(\mathbf{x}^{[i]}, y^{[i]}) | i = 1, \dots, n\} \quad \mathbf{x}^{[i]} \in \mathbb{R}^{m+1} \quad y^{[i]} \in \{0, 1\}$$

```
Result:  $\mathbf{w}$ 
 $\mathbf{w} = \mathbf{0}_{m+1}$ ;
for  $t = 1, \dots, T$  do
    for  $i = 1, \dots, n$  do
        instructions;
        if  $y^{[i]}(\mathbf{w} \cdot \mathbf{x}^{[i]}) \leq 0$  then
             $\mathbf{w} = \mathbf{w} + y^{[i]}\mathbf{x}^{[i]}$ 
        else
            do nothing
        end
    end
end
end
```

In [2]:

```
z = 0.  
for i in range(len(x)):  
    z += x[i] * w[i]  
  
print(z)
```

2.2

In [3]:

```
z = sum(x_i*w_i for x_i, w_i in zip(x, w))  
print(z)
```

2.2

In [4]:

```
import numpy as np  
  
x_vec, w_vec = np.array(x), np.array(w)  
  
z = (x_vec.transpose()).dot(w_vec)  
print(z)  
  
z = x_vec.dot(w_vec)  
print(z)
```

2.2

2.2

In [6]: `%timeit -r 100 -n 10 forloop(x, w)`

38.9 ms \pm 1.32 ms per loop (mean \pm std. dev. of 100 runs, 10 loops each)

In [7]: `%timeit -r 100 -n 10 listcomprehension(x, w)`

29.7 ms \pm 842 μ s per loop (mean \pm std. dev. of 100 runs, 10 loops each)

In [8]: `%timeit -r 100 -n 10 vectorized(x_vec, w_vec)`

46.8 μ s \pm 8.07 μ s per loop (mean \pm std. dev. of 100 runs, 10 loops each)

Implementación

```
3
4 class Perceptron():
5     def __init__(self, num_features):
6         self.num_features = num_features
7         self.weights = np.zeros((num_features, 1), dtype=np.float)
8         self.bias = np.zeros(1, dtype=np.float)
9
10
11
12
13
14
15
16
17
18     def forward(self, x):
19         linear = np.dot(x, self.weights) + self.bias
20         predictions = np.where(linear > 0., 1, 0)
21         return predictions
22
23     def backward(self, x, y):
24         predictions = self.forward(x)
25         errors = y - predictions
26         return errors
27
28     def train(self, x, y, epochs):
29         for e in range(epochs):
30
31             for i in range(y.shape[0]):
32
33                 errors = self.backward(x[i].reshape(1, self.num_features),
34 y[i]).reshape(-1)
35                 self.weights += (errors * x[i]).reshape(self.num_features,
36 1)
37                 self.bias += errors
38
39     def evaluate(self, x, y):
40         predictions = self.forward(x).reshape(-1)
41         accuracy = np.sum(predictions == y) / y.shape[0]
42         return accuracy
```

```
3
4 class Perceptron():
5     def __init__(self, num_features):
6         self.num_features = num_features
7
8
9         self.weights = torch.zeros(num_features, 1,
10                                     dtype=torch.float32, device=device)
11         self.bias = torch.zeros(1, dtype=torch.float32, device=device)
12
13         # placeholder vectors so they don't
14         # need to be recreated each time
15         self.ones = torch.ones(1)
16         self.zeros = torch.zeros(1)
17
18     def forward(self, x):
19         linear = torch.add(torch.mm(x, self.weights), self.bias)
20         predictions = torch.where(linear > 0., self.ones, self.zeros)
21         return predictions
22
23     def backward(self, x, y):
24         predictions = self.forward(x)
25         errors = y - predictions
26         return errors
27
28     def train(self, x, y, epochs):
29         for e in range(epochs):
30
31             for i in range(y.shape[0]):
32
33                 # use view because backward expects a matrix (i.e., 2D tensor)
34                 errors = self.backward(x[i].reshape(1, self.num_features),
35 y[i]).reshape(-1)
36                 self.weights += (errors * x[i]).reshape(self.num_features,
37 1)
38                 self.bias += errors
39
40     def evaluate(self, x, y):
41         predictions = self.forward(x).reshape(-1)
42         accuracy = torch.sum(predictions == y).float() / y.shape[0]
43         return accuracy
```

[...] Where a perceptron had been trained to distinguish between - this was for military purposes - it was looking at a scene of a forest in which there were camouflaged tanks in one picture and no camouflaged tanks in the other. And the perceptron - after a little training - made a 100% correct distinction between these two different sets of photographs. Then they were embarrassed a few hours later to discover that the two rolls of film had been developed differently. And so these pictures were just a little darker than all of these pictures and the perceptron was just measuring the total amount of light in the scene. But it was very clever of the perceptron to find some way of making the distinction.

-- Marvin Minsky, AI researcher & author of the "Perceptrons" book

Source: <https://www.webofstories.com/play/marvin.minsky/122>



<https://qph.fs.quoracdn.net/main-qimg-305eb8136c4a20f348bb7ab465bc2e10>



<http://theconversation.com/want-to-beat-climate-change-protect-our-natural-forests-121491>

Como principio general de aprendizaje, aplicable a todos los modelos comunes de neuronas y arquitecturas de redes neuronales (profundas y no profundas):

Sea

$$\mathcal{D} = ((\mathbf{x}^{[1]}, y^{[1]}), (\mathbf{x}^{[2]}, y^{[2]}), \dots, (\mathbf{x}^{[n]}, y^{[n]})) \in (\mathbb{R}^m \times \{0, 1\})^n$$

Modo *On-line*

Result: \mathbf{w}, b

$\mathbf{w} := \mathbf{0} \in \mathbb{R}^m, \mathbf{b} := 0;$

for $t = 1, \dots, T$ **do**

for $i = 1, \dots, n$ **do**

 cálculo de la salida;

 cálculo del error ;

 actualización de parámetros \mathbf{w}, b ;

end

end

Modo *On-line*

```
Result:  $\mathbf{w}, b$   
 $\mathbf{w} := \mathbf{0} \in \mathbb{R}^m, \mathbf{b} := 0;$   
for  $t = 1, \dots, T$  do  
    for  $i = 1, \dots, n$  do  
        cálculo de la salida;  
        cálculo del error ;  
        actualización de  
        parámetros  $\mathbf{w}, b$ ;  
    end  
end
```

Modo *On-line* II

```
Result:  $\mathbf{w}, b$   
 $\mathbf{w} := \mathbf{0} \in \mathbb{R}^m, \mathbf{b} := 0;$   
for  $j$  iteraciones do  
    Elija un  $(\mathbf{x}^{[i]}, y^{[i]}) \in \mathcal{D}$   
    aleatorio;  
    cálculo de la salida;  
    cálculo del error ;  
    actualización  $\mathbf{w}, b$ ;  
end
```

Modo *Batch*

Modo *On-line*

Result: \mathbf{w}, b

$\mathbf{w} := \mathbf{0} \in \mathbb{R}^m, \mathbf{b} := 0;$

for $t = 1, \dots, T$ **do**

for $i = 1, \dots, n$ **do**

 cálculo de la salida;

 cálculo del error ;

 actualización de
 parámetros \mathbf{w}, b ;

end

end

Result: \mathbf{w}, b

$\mathbf{w} := \mathbf{0} \in \mathbb{R}^m, \mathbf{b} := 0;$

for $t = 1, \dots, T$ **do**

$\Delta \mathbf{w} := 0, \Delta b := 0;$

for $i = 1, \dots, n$ **do**

 cálculo de la salida;

 cálculo del error ;

 actualización de
 parámetros $\Delta \mathbf{w}, \Delta b$;

end

 actualización de parámetros

\mathbf{w}, b ;

$\mathbf{w} := \mathbf{w} + \Delta \mathbf{w} \quad b := b + \Delta b$;

end

Modo *minibatch*

Modo más común en Deep Learning. Combina *On-line* y *Batch*.

$$\mathcal{D} = ((\mathbf{x}^{[1]}, y^{[1]}), (\mathbf{x}^{[2]}, y^{[2]}), \dots, (\mathbf{x}^{[n]}, y^{[n]})) \in (\mathbb{R}^m \times \{0, 1\})^n$$

Result: \mathbf{w}, b

$\mathbf{w} := \mathbf{0} \in \mathbb{R}^m, \mathbf{b} := 0;$

for $t = 1, \dots, T$ **do**

for $j = 1, \dots, n/k$ **do**

$\Delta \mathbf{w} := 0, \Delta b := 0;$

for $\{(\mathbf{x}^{[i]}, y^{[i]}), \dots, (\mathbf{x}^{[i+k]}, y^{[i+k]})\} \subset D$ **do**

 cálculo de la salida;

 cálculo del error ;

 actualización de $\Delta \mathbf{w}, \Delta b;$

end

 actualización $\mathbf{w}, b;$

$\mathbf{w} := \mathbf{w} + \Delta \mathbf{w};$

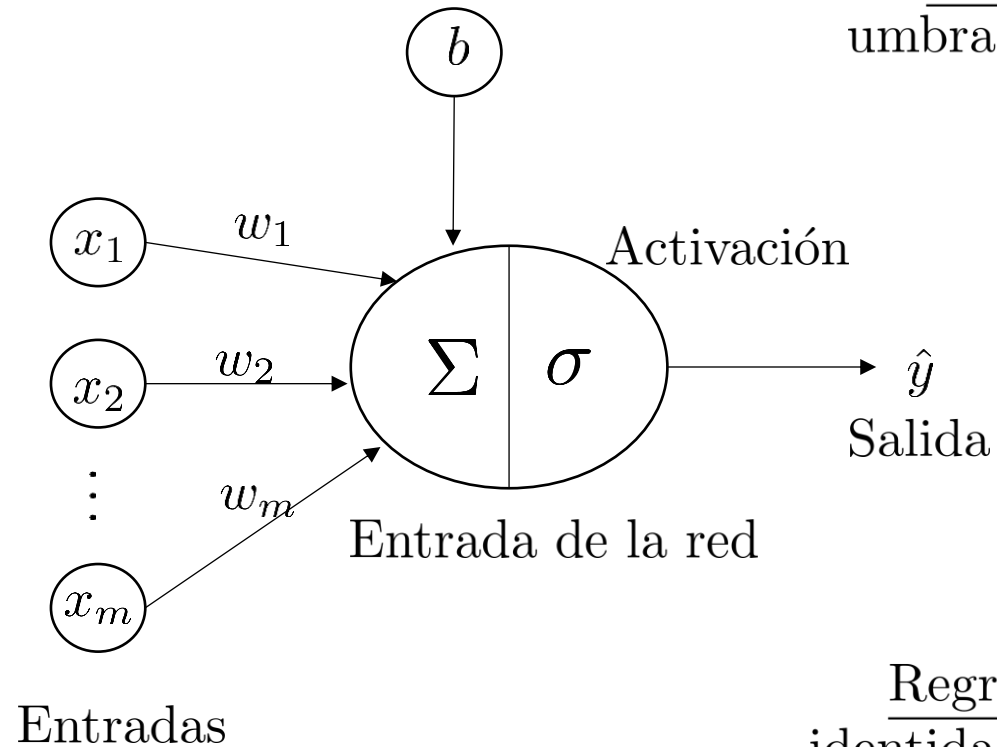
$b := b + \Delta w;$

end

end

Regresión lineal

Perceptrón: la función de activación es la función de umbral. La salida es la etiqueta binaria $\hat{y} \in \{0, 1\}$



Regresión lineal: la función de activación es la función identidad $\sigma(x) = x$. La salida es el número $\hat{y} \in \mathbb{R}$

Regresión Lineal

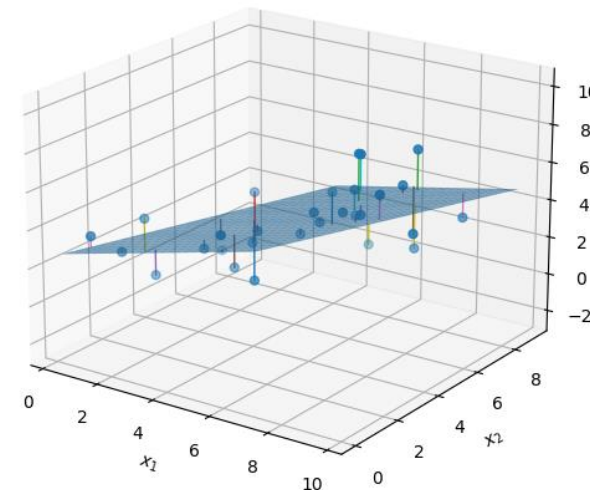
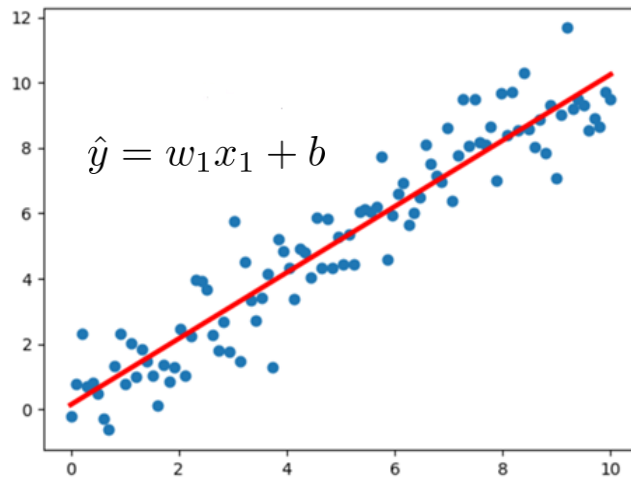
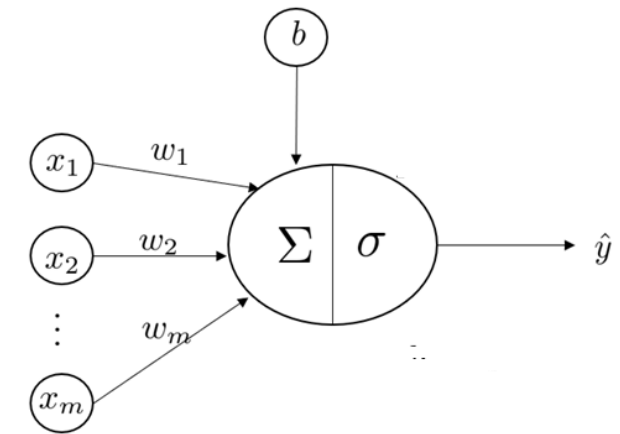
$$\hat{y} = w_1 x_1 + b$$

Describe la relación entre una variable dependiente \hat{y} y una única variable independiente x .

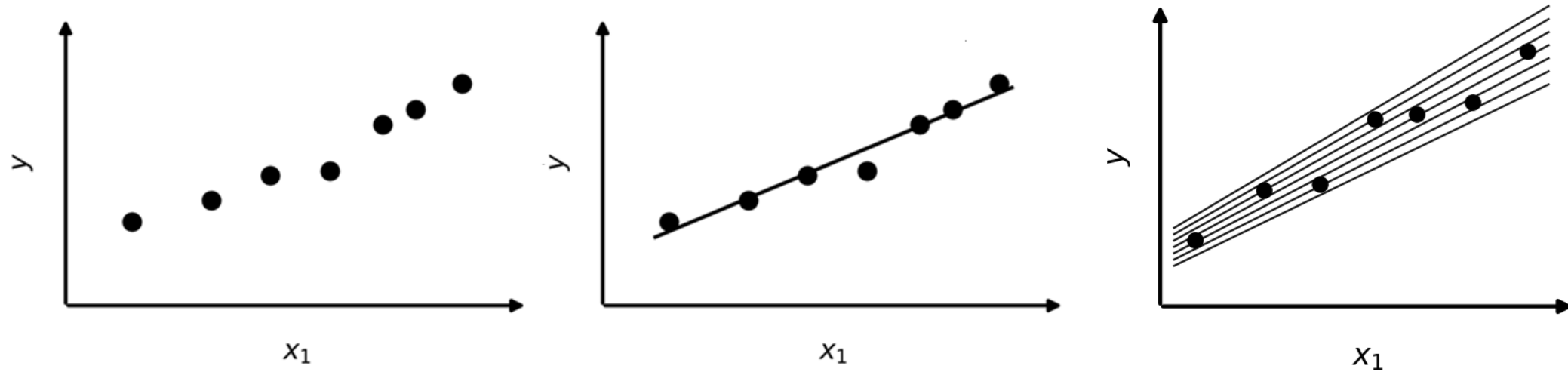
Regresión Lineal Múltiple

$$\hat{y} = b + w_1 x_1 + w_2 x_2 + \dots + w_m x_m$$

Relaciona una variable dependiente \hat{y} con dos o más variables independientes.



Regresión lineal

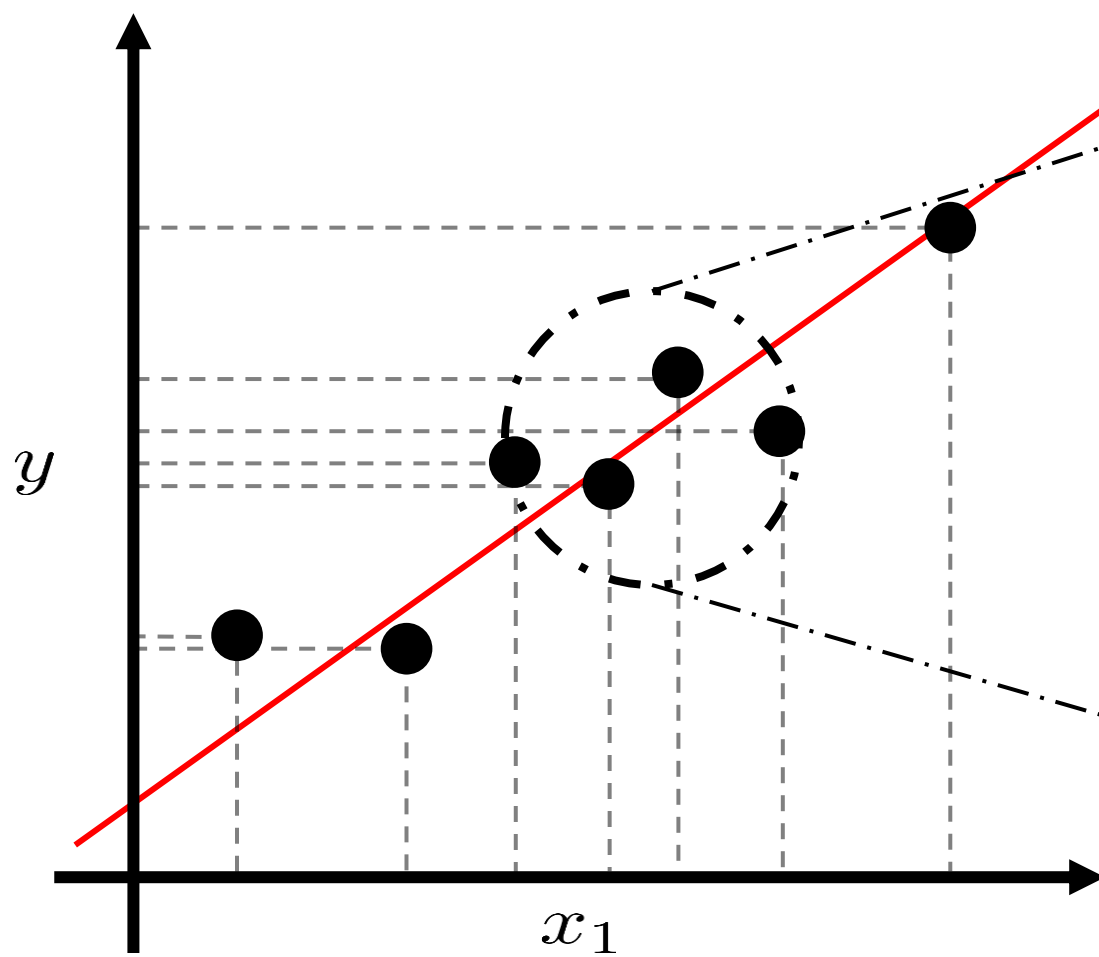


¿Cuál es la línea que mejor se ajusta a los datos?

Minimizar la distancia entre la línea y cada punto

Diferencia entre el valor predicho y el valor real

Regresión lineal

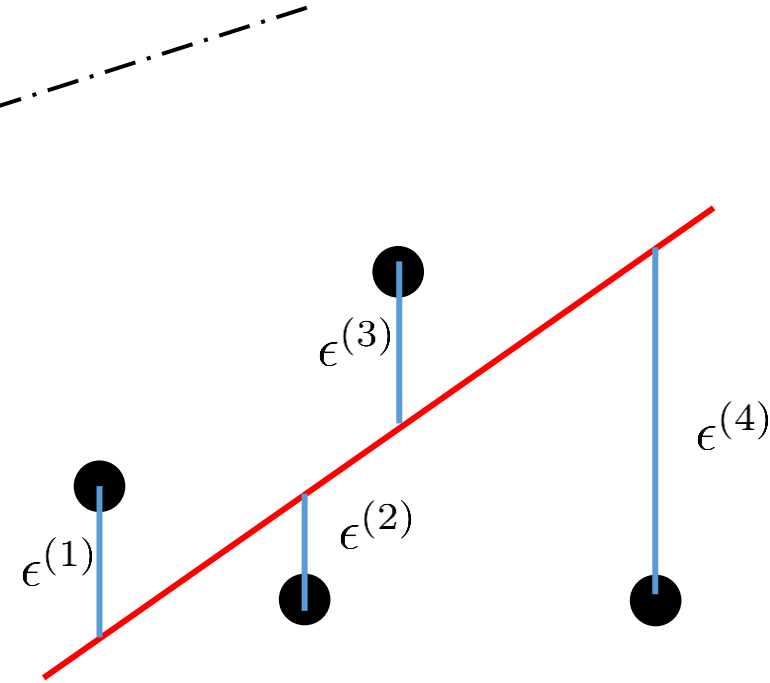


$$\hat{y} = w_1 x_1 + b$$

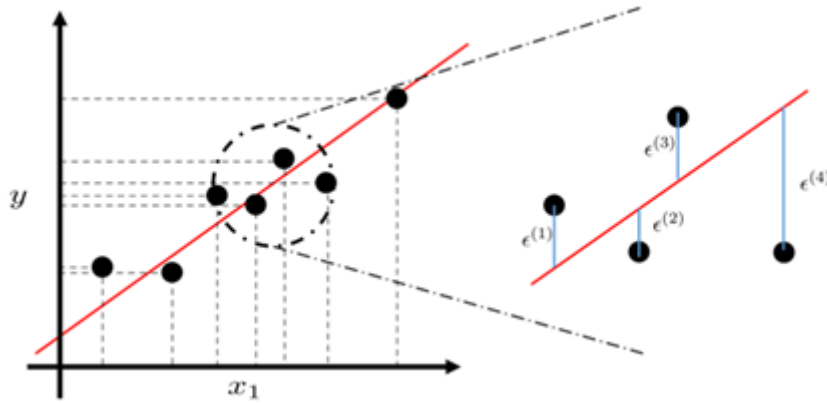


$$\epsilon^{(i)} = y^{(i)} - \hat{y}^{(i)}$$

$$\epsilon^{(i)} = y^{(i)} - (w_1 x_1^{(i)} + b)$$



Regresión lineal



Minimizar la distancia entre la línea y cada punto

Función de costo:

$$L(w_1, b) = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - (w_1 x^{(i)} + b))^2$$

$$L(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - (\mathbf{w}^\top \mathbf{x}^{(i)} + b))^2$$

Regresión lineal (mínimos cuadrados)

El modelo de regresión lineal que se puede usar es el siguiente:

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y},$$

asumiendo que el bias está incluido en \mathbf{w} , y la matriz de diseño tiene un vector adicional. Es posible resolver este modelo usando *ecuaciones normales*, i.e., solución analítica.

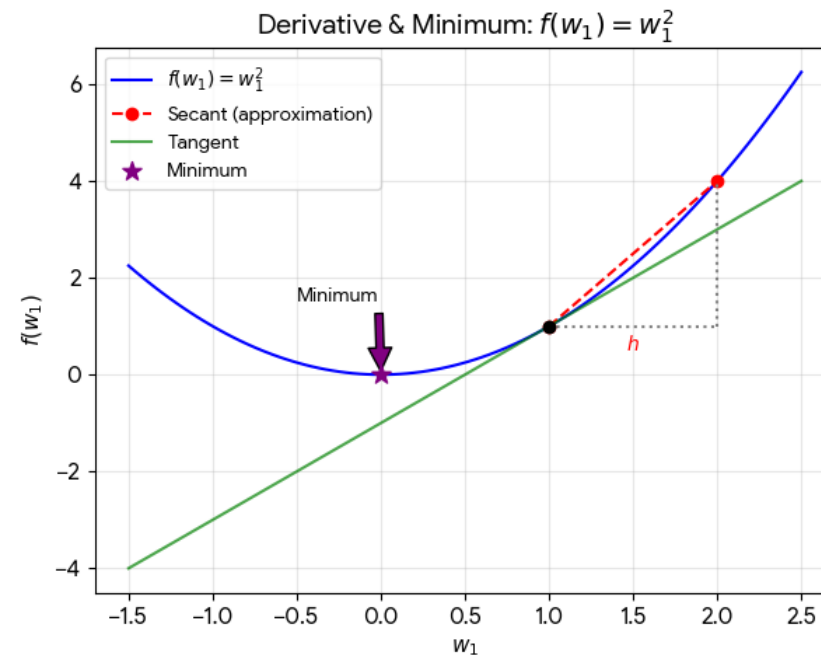
Regresión lineal (mínimos cuadrados)

Una mejor manera

Una mejor forma de hacerlo, es analizando qué efecto tiene un cambio en un parámetro en el desempeño predictivo (pérdida) del modelo; posteriormente, se cambiará solo un poco el peso en la dirección que mejora el desempeño (minimiza la pérdida). Esto se hace con bastantes pasos pequeños, hasta que la pérdida no se decremente más.

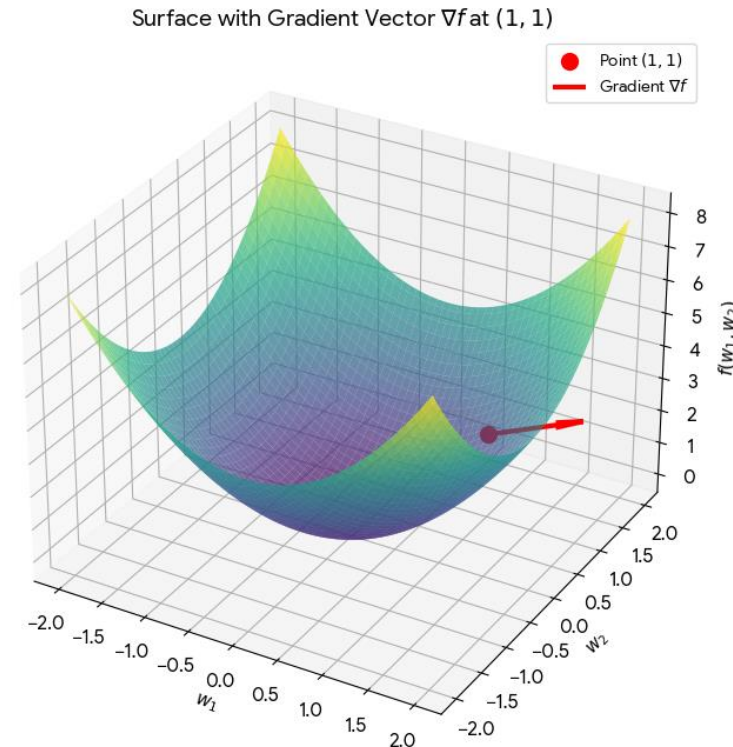
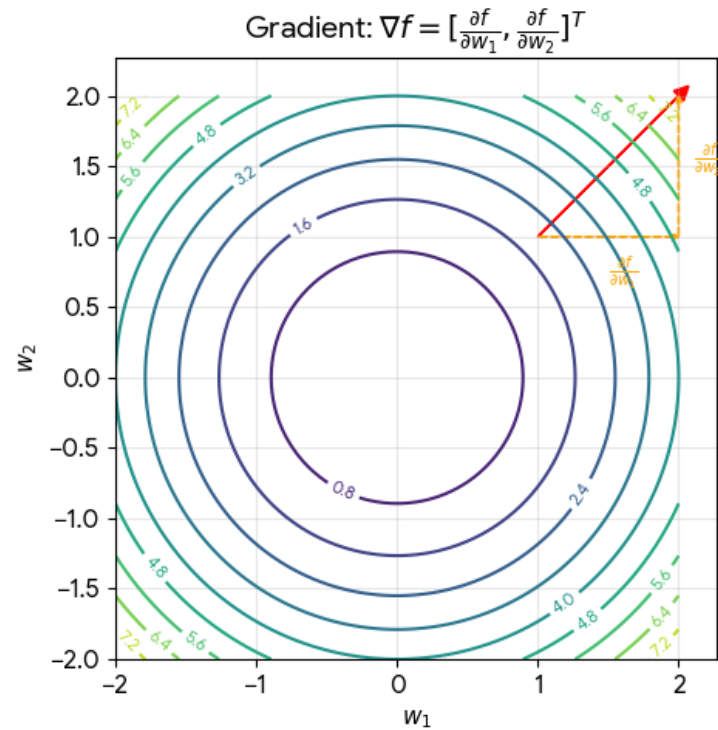
Derivadas de funciones univariadas

$$\frac{d}{dw_1} f(w_1) = \lim_{h \rightarrow 0} \frac{f(w_1+h) - f(w_1)}{h}$$



Gradiente: Derivadas de funciones multivariadas

$$\nabla f(w_1, w_2) = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \end{bmatrix} = \begin{bmatrix} \lim_{h \rightarrow 0} \frac{f(w_1+h, w_2) - f(w_1, w_2)}{h} \\ \lim_{h \rightarrow 0} \frac{f(w_1, w_2+h) - f(w_1, w_2)}{h} \end{bmatrix}$$



Gradiente: Derivadas de funciones multivariantes

Si se tiene la función $f(x, y, z, \dots)$,

el gradiente tendrá la forma $\Delta f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \\ \frac{\partial f}{\partial z} \\ \vdots \end{bmatrix}$.

Ejemplo

Para la función $f(x, y) = x^2y + y$ tendremos el gradiente $\Delta f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}$,
donde

$$\frac{\partial f}{\partial x} = \frac{\partial}{\partial x} x^2y + y = 2xy$$

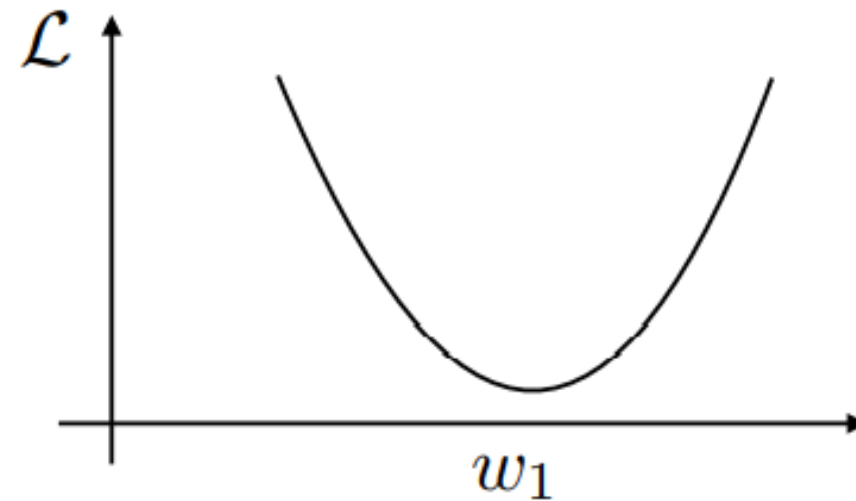
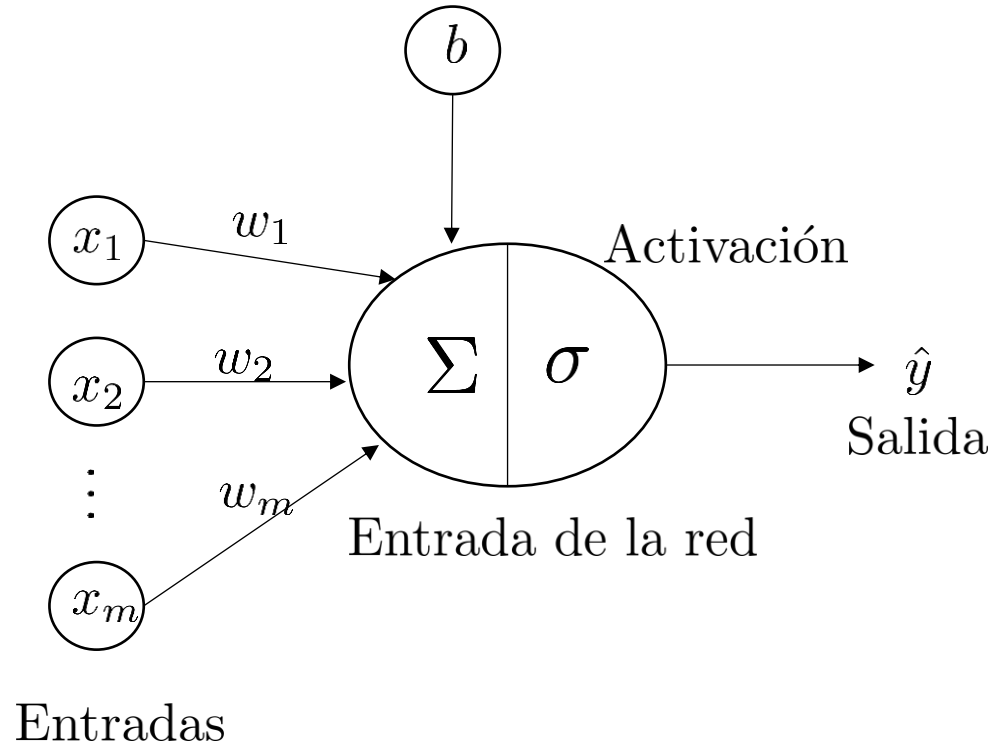
y

$$\frac{\partial f}{\partial y} = \frac{\partial}{\partial y} x^2y + y = x^2 + 1.$$

Por tanto, el gradiente de la función f estará definido como

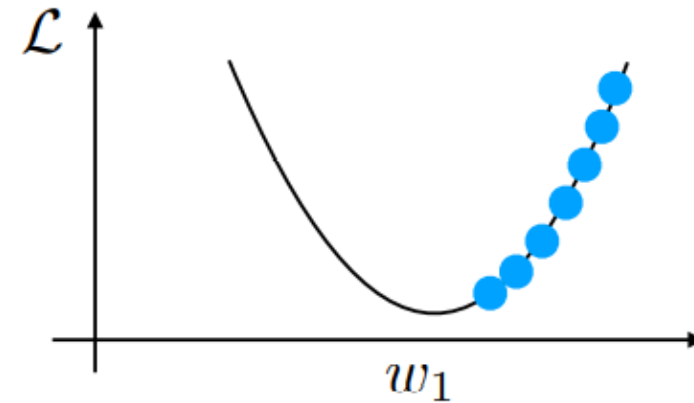
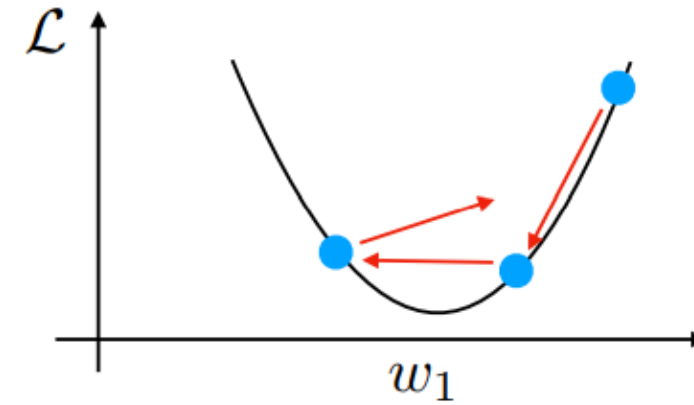
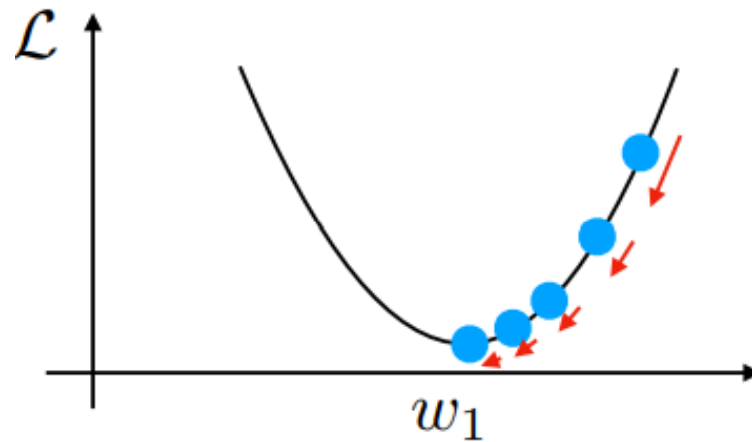
$$\Delta f(x, y) = \begin{bmatrix} 2xy \\ x^2 + 1 \end{bmatrix}$$

Entrenamiento del regresor lineal con gradiente descendente



Función de pérdida convexa
$$\mathcal{L}(\mathbf{w}, b) = \sum_i (\hat{y}^{[i]} - y^{[i]})^2$$

Entrenamiento del regresor lineal con gradiente descendente



Regla de aprendizaje del perceptrón

```
Result:  $w, b$ 
 $w := \mathbf{0} \in \mathbb{R}^m, b := 0;$ 
for  $t = 1, \dots, T$  do
  for  $i = 1, \dots, n$  do
     $\hat{y}^{[i]} := \sigma(\mathbf{x}^{[i]\top} \mathbf{w} + b);$ 
     $err := (y^{[i]} - \hat{y}^{[i]});$ 
     $\mathbf{w} := \mathbf{w} + err \times \mathbf{x}^{[i]},$ 
     $b := b + err;$ 
  end
end
end
```

Gradiente descendente estocástico

```
Result:  $w, b$ 
 $\mathbf{w} := \mathbf{0} \in \mathbb{R}^m, b := 0;$ 
for  $t = 1, \dots, T$  do
  for  $i = 1, \dots, n$  do
     $\hat{y}^{[i]} := \sigma(\mathbf{x}^{[i]\top} \mathbf{w} + b);$ 
     $\Delta_{\mathbf{w}} \mathcal{L} = (y^{[i]} - \hat{y}^{[i]}) \mathbf{x}^{[i]} ;$ 
     $\Delta_b \mathcal{L} = (y^{[i]} - \hat{y}^{[i]});$ 
     $\mathbf{w} := \mathbf{w} + \eta \times (-\Delta_{\mathbf{w}} \mathcal{L}) ;$  gradiente negativo
     $b := b + \eta \times (-\Delta_b \mathcal{L});$ 
  end
end
end
```

tasa de aprendizaje

Vectorizado

Result: \mathbf{w}, b

$\mathbf{w} := \mathbf{0} \in \mathbb{R}^m, \mathbf{b} := 0;$

for $t = 1, \dots, T$ **do**

for $i = 1, \dots, n$ **do**

$\hat{y}^{[i]} := \sigma(\mathbf{x}^{[i]\top} \mathbf{w} + b);$

$\Delta_{\mathbf{w}} \mathcal{L} = (y^{[i]} - \hat{y}^{[i]}) \mathbf{x}^{[i]} ;$

$\Delta_b \mathcal{L} = (y^{[i]} - \hat{y}^{[i]});$

$\mathbf{w} := \mathbf{w} + \eta \times (-\Delta_{\mathbf{w}} \mathcal{L}) ;$

$b := b + \eta \times (-\Delta_b \mathcal{L});$

end

end

Ciclo *for*

Result: \mathbf{w}, b

$\mathbf{w} := \mathbf{0} \in \mathbb{R}^m, \mathbf{b} := 0;$

for $t = 1, \dots, T$ **do**

$\Delta \mathbf{w} := \mathbf{0}, \Delta b := 0;$

for $i = 1, \dots, n$ **do**

$\hat{y}^{[i]} := \sigma(\mathbf{x}^{[i]\top} \mathbf{w} + b);$

end

for $j = 1, \dots, m$ **do**

$\frac{\partial \mathcal{L}}{\partial w_j} = (y^{[i]} - \hat{y}^{[i]}) x_j^{[i]};$

$w_j := w_j + \eta \times \left(-\frac{\partial \mathcal{L}}{\partial w_j}\right);$

end

$\frac{\partial \mathcal{L}}{\partial b} = (y^{[i]} - \hat{y}^{[i]});$

$b := b + \eta \times \left(-\frac{\partial \mathcal{L}}{\partial b}\right);$

end

Modo *On-line*

```
Result:  $\mathbf{w}, b$   
 $\mathbf{w} := \mathbf{0} \in \mathbb{R}^m, \mathbf{b} := 0;$   
for  $t = 1, \dots, T$  do  
   $\Delta \mathbf{w} := 0, \Delta b := 0;$   
  for  $i = 1, \dots, n$  do  
     $\hat{y}^{[i]} := \sigma(\mathbf{x}^{[i]\top} \mathbf{w} + b);$   
  end  
  for  $j = 1, \dots, m$  do  
     $\frac{\partial \mathcal{L}}{\partial w_j} = (y^{[i]} - \hat{y}^{[i]})x_j^{[i]};$   
     $w_j := w_j + \eta \times \left(-\frac{\partial \mathcal{L}}{\partial w_j}\right);$   
  end  
   $\frac{\partial \mathcal{L}}{\partial b} = (y^{[i]} - \hat{y}^{[i]});$   
   $b := b + \eta \times \left(-\frac{\partial \mathcal{L}}{\partial b}\right);$   
end
```

Coincidentalmente, parece casi el mismo de la regla del perceptrón, exceptuando que la predicción es un número real y se tiene una tasa de aprendizaje. Esta regla de aprendizaje se llama **gradiente descendente estocástico**

Dada la función de pérdida

$$\mathcal{L}(\mathbf{w}, b) = \sum_i (\hat{y}^{[i]} - y^{[i]})^2,$$

la derivada estaría dada por

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_j} &= \frac{\partial}{\partial w_j} \sum_i (\hat{y}^{[i]} - y^{[i]})^2 \\ &= \frac{\partial}{\partial w_j} \sum_i (\sigma(\mathbf{w}^\top \mathbf{x}^{[i]}) - y^{[i]})^2 \\ &= \sum_i 2(\sigma(\mathbf{w}^\top \mathbf{x}^{[i]}) - y^{[i]}) \frac{\partial}{\partial w_j} (\sigma(\mathbf{w}^\top \mathbf{x}^{[i]}) - y^{[i]}) \\ &= \sum_i 2(\sigma(\mathbf{w}^\top \mathbf{x}^{[i]}) - y^{[i]}) \frac{d\sigma}{d(\mathbf{w}^\top \mathbf{x}^{[i]})} \frac{\partial}{\partial w_j} \mathbf{w}^\top \mathbf{x}^{[i]} \\ &= \sum_i 2(\sigma(\mathbf{w}^\top \mathbf{x}^{[i]}) - y^{[i]}) \frac{d\sigma}{d(\mathbf{w}^\top \mathbf{x}^{[i]})} x_j^{[i]} \\ &= \sum_i 2(\sigma(\mathbf{w}^\top \mathbf{x}^{[i]}) - y^{[i]}) x_j^{[i]} \end{aligned}$$

Note que la función de activación para el caso de la regresión lineal es la función identidad.

Derivada de la función de pérdida para regresión lineal: caso alternativo

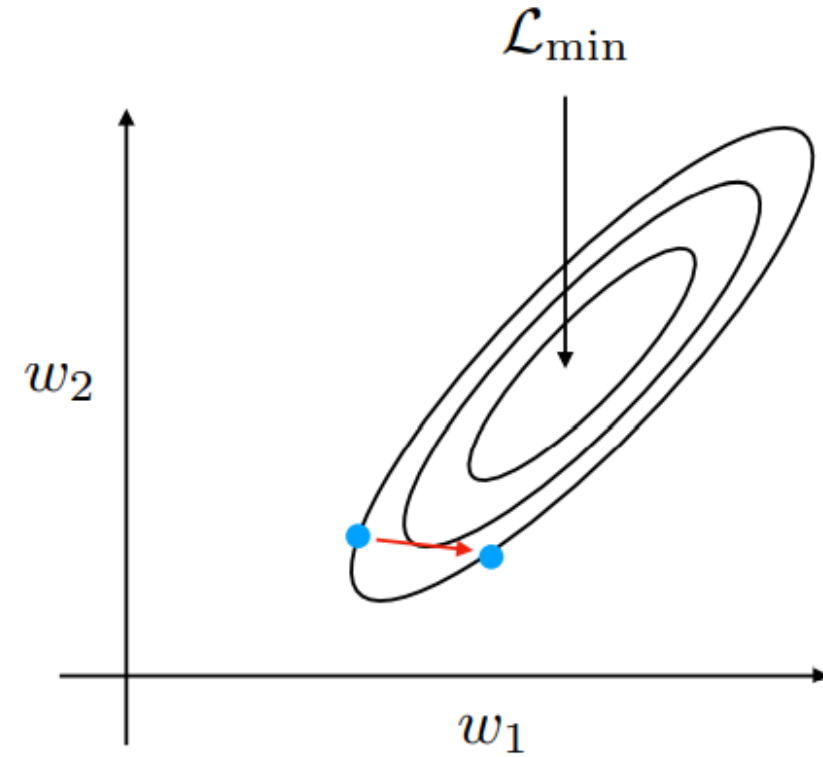
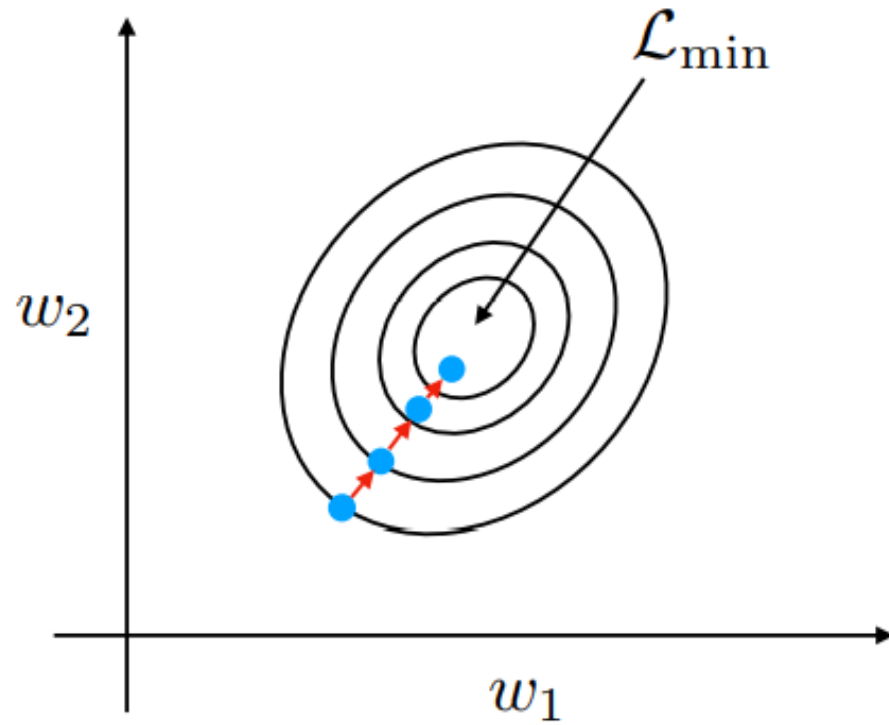
A menudo, la función de pérdida se escala por un factor de $1/2$ por conveniencia:

$$\mathcal{L}(\mathbf{w}, b) = \frac{1}{2n} \sum_i (\hat{y}^{[i]} - y^{[i]})^2.$$

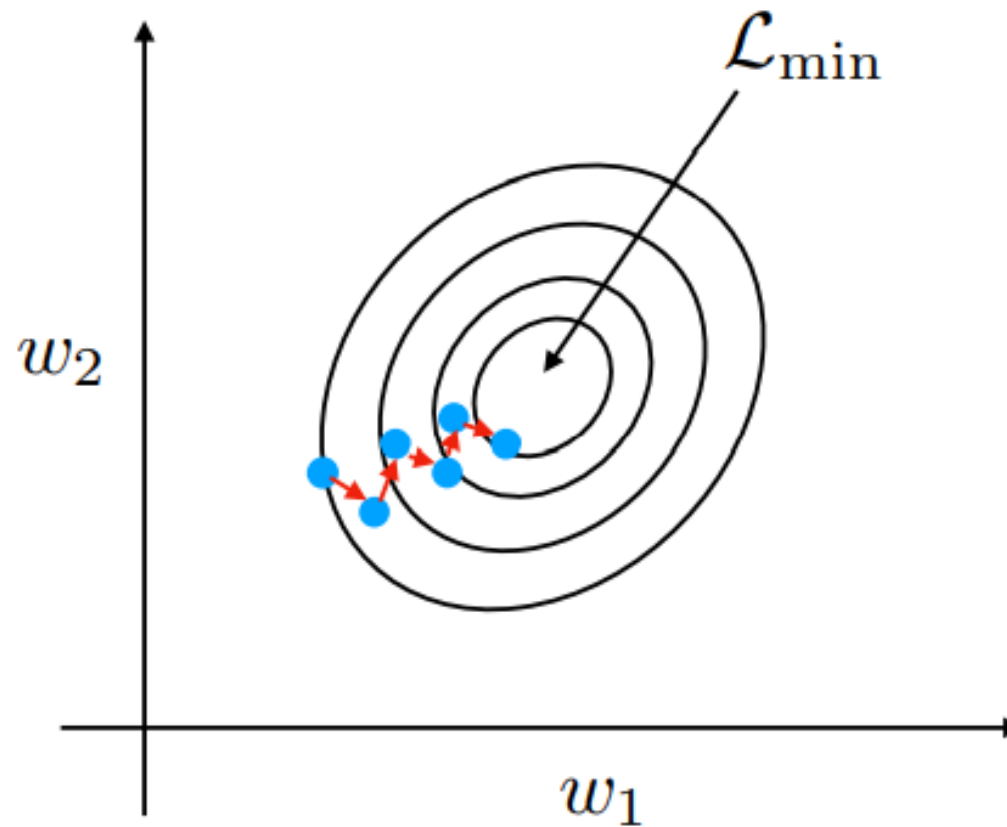
La derivada estaría dada por

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_j} &= \frac{\partial}{\partial w_j} \frac{1}{2n} \sum_i (\hat{y}^{[i]} - y^{[i]})^2 \\ &= \frac{\partial}{\partial w_j} \sum_i \frac{1}{2n} (\sigma(\mathbf{w}^\top \mathbf{x}^{[i]}) - y^{[i]})^2 \\ &= \sum_i \frac{1}{n} (\sigma(\mathbf{w}^\top \mathbf{x}^{[i]}) - y^{[i]}) \frac{\partial}{\partial w_j} (\sigma(\mathbf{w}^\top \mathbf{x}^{[i]}) - y^{[i]}) \\ &= \sum_i \frac{1}{n} (\sigma(\mathbf{w}^\top \mathbf{x}^{[i]}) - y^{[i]}) \frac{d\sigma}{d(\mathbf{w}^\top \mathbf{x}^{[i]})} \frac{\partial}{\partial w_j} \mathbf{w}^\top \mathbf{x}^{[i]} \\ &= \sum_i \frac{1}{n} (\sigma(\mathbf{w}^\top \mathbf{x}^{[i]}) - y^{[i]}) \frac{d\sigma}{d(\mathbf{w}^\top \mathbf{x}^{[i]})} x_j^{[i]} \\ &= \sum_i \frac{1}{n} (\sigma(\mathbf{w}^\top \mathbf{x}^{[i]}) - y^{[i]}) x_j^{[i]} \end{aligned}$$

Gradiente descendente por lotes



Gradiente descendente estocástico



Entrenamiento de una red neuronal de una sola capa con gradiente descendente

Widrow and Hoff's ADALINE (1960)

A nicely differentiable neuron model

Widrow, B., & Hoff, M. E. (1960). *Adaptive switching circuits* (No. TR-1553-1). Stanford Univ Ca Stanford Electronics Labs.

Widrow, B. (1960). *Adaptive" adaline" Neuron Using Chemical" memistors."*.

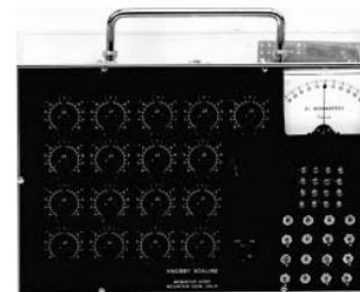
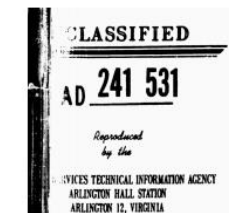


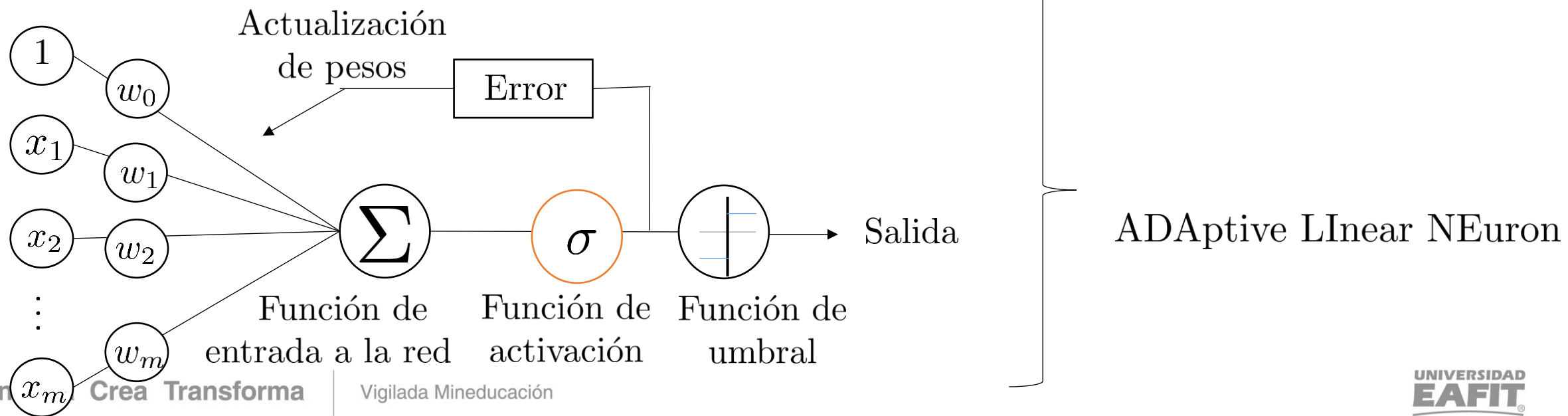
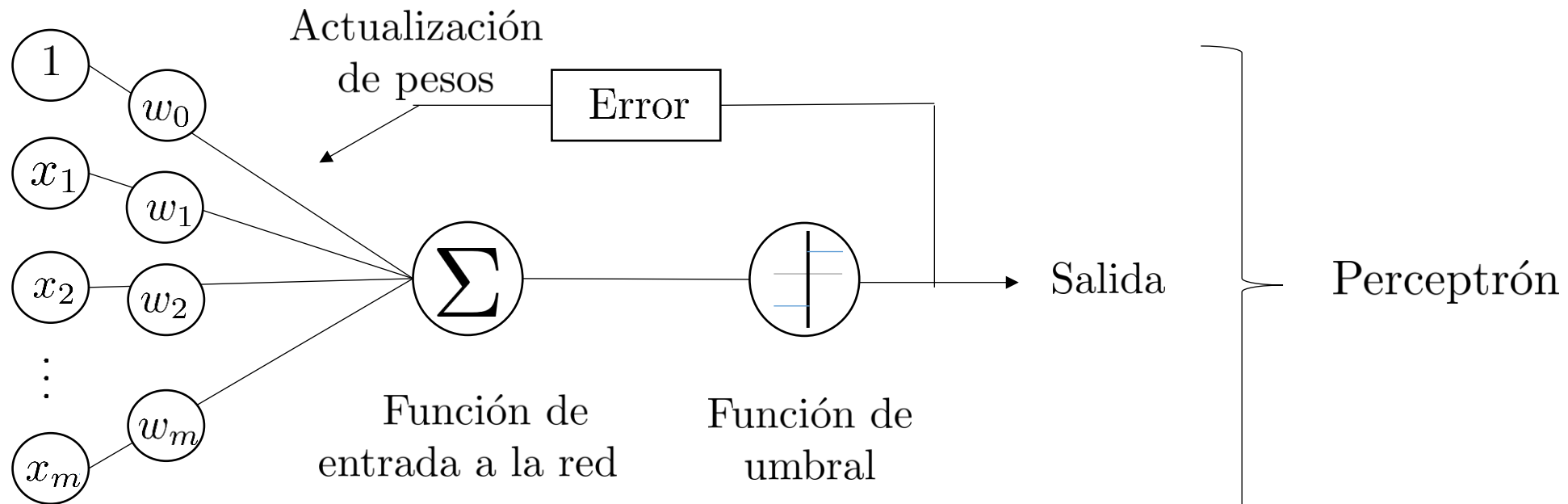
Image source: https://www.researchgate.net/profile/Alexander_Magoun2/publication/265789430/figure/fig2/AS:630233355178778001470551421849/ADALINE-An-adaptive-linear-neuron-Manually-adapted-synapses-Designed-and-built-by-Ted.png



THIS REPORT HAS BEEN DELIMITED
AND CLEARED FOR PUBLIC RELEASE
UNDER DOD DIRECTIVE 5200.20 AND
NO RESTRICTIONS ARE IMPOSED UPON
ITS USE AND DISCLOSURE.

DISTRIBUTION STATEMENT A

APPROVED FOR PUBLIC RELEASE;
DISTRIBUTION UNLIMITED.



Algunos tutoriales sobre PyTorch

- <https://pytorch.org/tutorials/>
- https://pytorch.org/tutorials/beginner/deep_learning_60min_blitz.html
- <https://pytorch.org/assets/deep-learning/Deep-Learning-with-PyTorch.pdf>

y el foro de discusión

<https://discuss.pytorch.org/>

