



# IIC 2433 Minería de Datos

<https://github.com/marcelomendoza/IIC2433>

- GRADIENTE DESCENDENTE -

## Regresión logística

Si el ajuste es adecuado:

$$P(y \mid \mathbf{x}) = \theta(y \cdot \mathbf{w}^T \mathbf{x})$$

Luego, podemos expresar la función de verosimilitud:

$$P(y_1, \dots, y_N \mid \mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{n=1}^N P(y_n \mid \mathbf{x}_n).$$

# Regresión logística

Maximizamos la función de verosimilitud:

$$\max \quad \prod_{n=1}^N P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \max \quad \ln \left( \prod_{n=1}^N P(y_n \mid \mathbf{x}_n) \right)$$

$$\equiv \max \quad \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

# Regresión logística

Maximizamos la función de verosimilitud:

$$\max \quad \prod_{n=1}^N P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \max \quad \ln \left( \prod_{n=1}^N P(y_n \mid \mathbf{x}_n) \right)$$

$$\equiv \max \quad \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \text{min} \quad - \frac{1}{N} \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

$$\equiv \min \quad \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{P(y_n \mid \mathbf{x}_n)}$$

$$\equiv \min \quad \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{\theta(y_n \cdot \mathbf{w}^T \mathbf{x}_n)}$$

# Regresión logística

Maximizamos la función de verosimilitud:

$$\max \quad \prod_{n=1}^N P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \max \quad \ln \left( \prod_{n=1}^N P(y_n \mid \mathbf{x}_n) \right)$$

$$\equiv \max \quad \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

$$\Leftrightarrow \text{min} \quad - \frac{1}{N} \sum_{n=1}^N \ln P(y_n \mid \mathbf{x}_n)$$

$$\equiv \min \quad \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{P(y_n \mid \mathbf{x}_n)}$$

$$\equiv \min \quad \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{\theta(y_n \cdot \mathbf{w}^T \mathbf{x}_n)}$$

$$\equiv \min \quad \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n \cdot \mathbf{w}^T \mathbf{x}_n})$$

$$\theta(s) = \frac{1}{1 + e^{-s}}.$$



# Regresión logística

Tenemos una expresión para:

Parámetros del modelo

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n \cdot \mathbf{w}^T \mathbf{x}_n}) \quad \text{Cross-entropy}$$

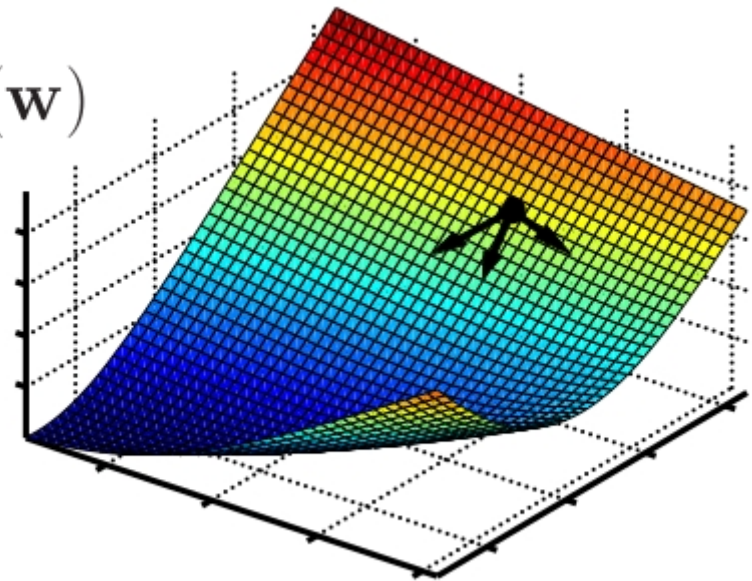
La función es convexa, por lo que podemos optimizarla de forma iterativa:

Idea del gradiente descendente:

$E_{\text{in}}(\mathbf{w})$

$$\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \hat{\mathbf{v}}$$

¿En cuál dirección conviene moverse?



## Gradiente descendente

$$\begin{aligned}\Delta E_{\text{in}} &= E_{\text{in}}(\mathbf{w}(t+1)) - E_{\text{in}}(\mathbf{w}(t)) \\ &= E_{\text{in}}(\mathbf{w}(t) + \eta \hat{\mathbf{v}}) - E_{\text{in}}(\mathbf{w}(t)) \\ &= \eta \nabla E_{\text{in}}(\mathbf{w}(t))^{\text{T}} \hat{\mathbf{v}} + O(\eta^2)\end{aligned}$$

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots + \frac{f^{(k+1)}(c)}{(k+1)!}(x - x_0)^{k+1}.$$

$$\begin{aligned}\Delta E_{\text{in}} &= E_{\text{in}}(\mathbf{w}(t+1)) - E_{\text{in}}(\mathbf{w}(t)) \\ &\quad \underbrace{f(x_0 + \eta \cdot x) - f(x_0)}_{\text{(expansión de Taylor de 1}^{\text{er}} \text{ orden)}} \\ &= f(x_0) + \nabla f(x_0) \cdot (x_0 + \eta \cdot x - x_0) + \dots\end{aligned}$$



## Gradiente descendente

$$\begin{aligned}
 \Delta E_{\text{in}} &= E_{\text{in}}(\mathbf{w}(t+1)) - E_{\text{in}}(\mathbf{w}(t)) \\
 &= E_{\text{in}}(\mathbf{w}(t) + \eta \hat{\mathbf{v}}) - E_{\text{in}}(\mathbf{w}(t)) \\
 &= \eta \nabla E_{\text{in}}(\mathbf{w}(t))^{\text{T}} \hat{\mathbf{v}} + O(\eta^2)
 \end{aligned}$$

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots + \frac{f^{(k+1)}(c)}{(k+1)!}(x - x_0)^{k+1}.$$

$$\Delta E_{\text{in}} = E_{\text{in}}(\mathbf{w}(t+1)) - E_{\text{in}}(\mathbf{w}(t))$$

$$\underbrace{f(x_0 + \eta \cdot x) - f(x_0)}$$

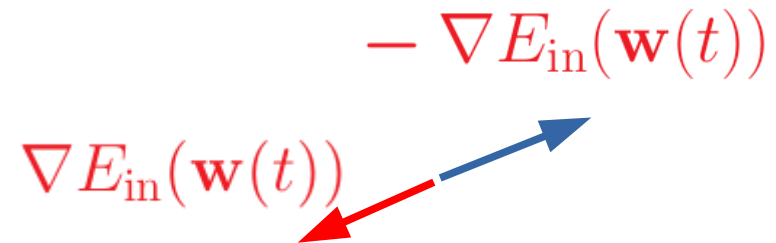
(expansión de Taylor de 1<sup>er</sup> orden)

$$\cancel{f(x_0)} + \nabla f(x_0) \cdot (\cancel{x_0} + \eta \cdot x - \cancel{x_0}) + \dots$$

## Gradiente descendente

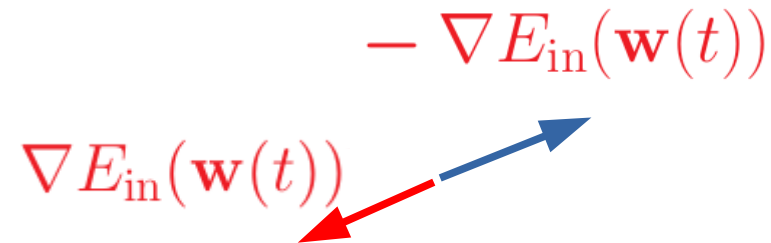
$$= \eta \underbrace{\nabla E_{\text{in}}(\mathbf{w}(t))^T}_{\text{Minimizado en}} \hat{\mathbf{v}} + O(\eta^2)$$

$$\text{Minimizado en } \hat{\mathbf{v}} = - \frac{\nabla E_{\text{in}}(\mathbf{w}(t))}{\|\nabla E_{\text{in}}(\mathbf{w}(t))\|}$$



## Gradiente descendente

$$= \eta \underbrace{\nabla E_{\text{in}}(\mathbf{w}(t))^T \hat{\mathbf{v}}}_{\text{Minimizado en } \hat{\mathbf{v}} = -\frac{\nabla E_{\text{in}}(\mathbf{w}(t))}{\|\nabla E_{\text{in}}(\mathbf{w}(t))\|}} + O(\eta^2)$$



Minimizado en  $\hat{\mathbf{v}} = -\frac{\nabla E_{\text{in}}(\mathbf{w}(t))}{\|\nabla E_{\text{in}}(\mathbf{w}(t))\|}$

Gradiente descendente

1: Initialize at step  $t = 0$  to  $\mathbf{w}(0)$ .

2: **for**  $t = 0, 1, 2, \dots$  **do**

3:   Compute the gradient

$$\mathbf{g}_t = \nabla E_{\text{in}}(\mathbf{w}(t)).$$

4:   Move in the direction  $\mathbf{v}_t = -\mathbf{g}_t$ .

5:   Update the weights:

$$\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \mathbf{v}_t.$$

6:   Iterate ‘until it is time to stop’.

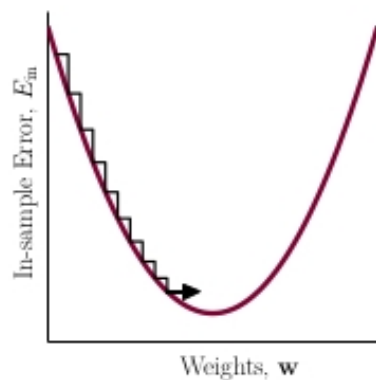
7: **end for**

8: Return the final weights.

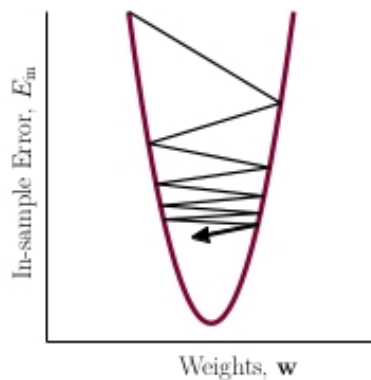
# Gradiente descendente

El efecto del **learning rate**:

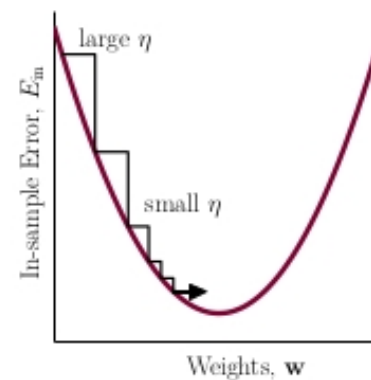
$\eta$  muy pequeño



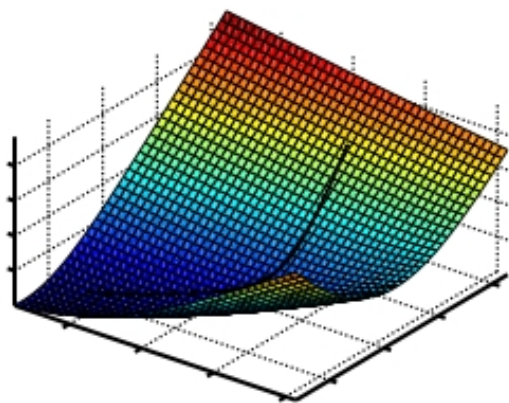
$\eta$  muy grande



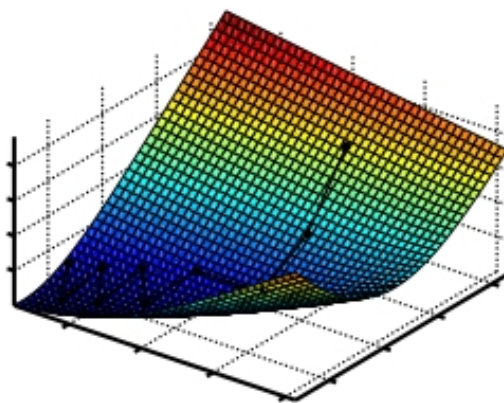
$\eta_t$



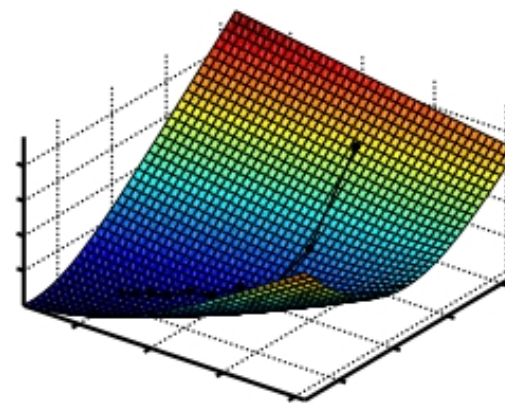
Adam



$\eta = 0.1$ ; 75 steps



$\eta = 2$ ; 10 steps



variable  $\eta_t$ ; 10 steps

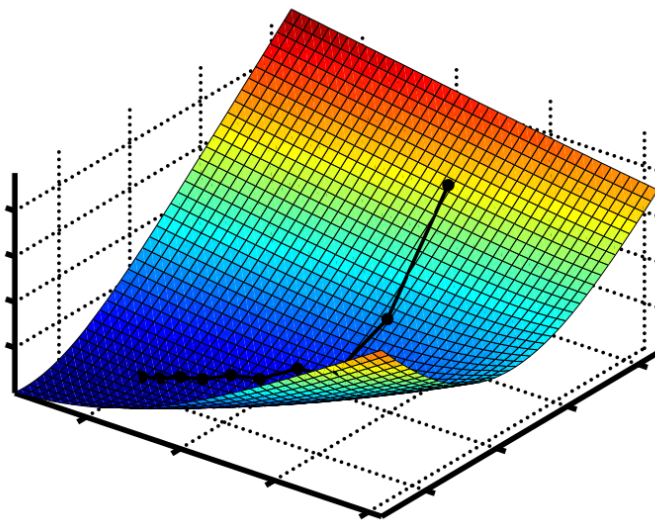
# Gradiente descendente

Una variación de gradiente descendente considera la evaluación del gradiente en una muestra del training set.

$$\mathbf{w}(t+1) \leftarrow \mathbf{w}(t) - \eta \nabla_{\mathbf{w}} e(\mathbf{w}, \mathbf{x}_*, y_*)$$

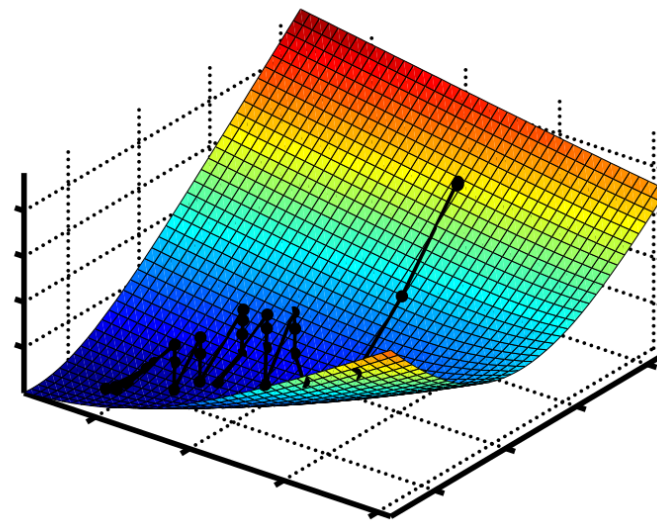
Dado que la muestra se toma al azar, se le denomina gradiente descendente estocástico.

GD



$\eta = 6$   
10 steps  
 $N = 10$

SGD



$\eta = 2$   
30 steps

Puede ayudar a escapar  
de óptimos locales

## Regresión logística + gradiente descendente

Para regresión logística encontramos que:

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n \cdot \mathbf{w}^T \mathbf{x}_n})$$

Muestre que:

$$\begin{aligned} \nabla E_{\text{in}}(\mathbf{w}) &= -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T \mathbf{x}_n}} \\ &= \frac{1}{N} \sum_{n=1}^N -y_n \mathbf{x}_n \theta(-y_n \mathbf{w}^T \mathbf{x}_n). \end{aligned}$$

## Regresión logística + gradiente descendente

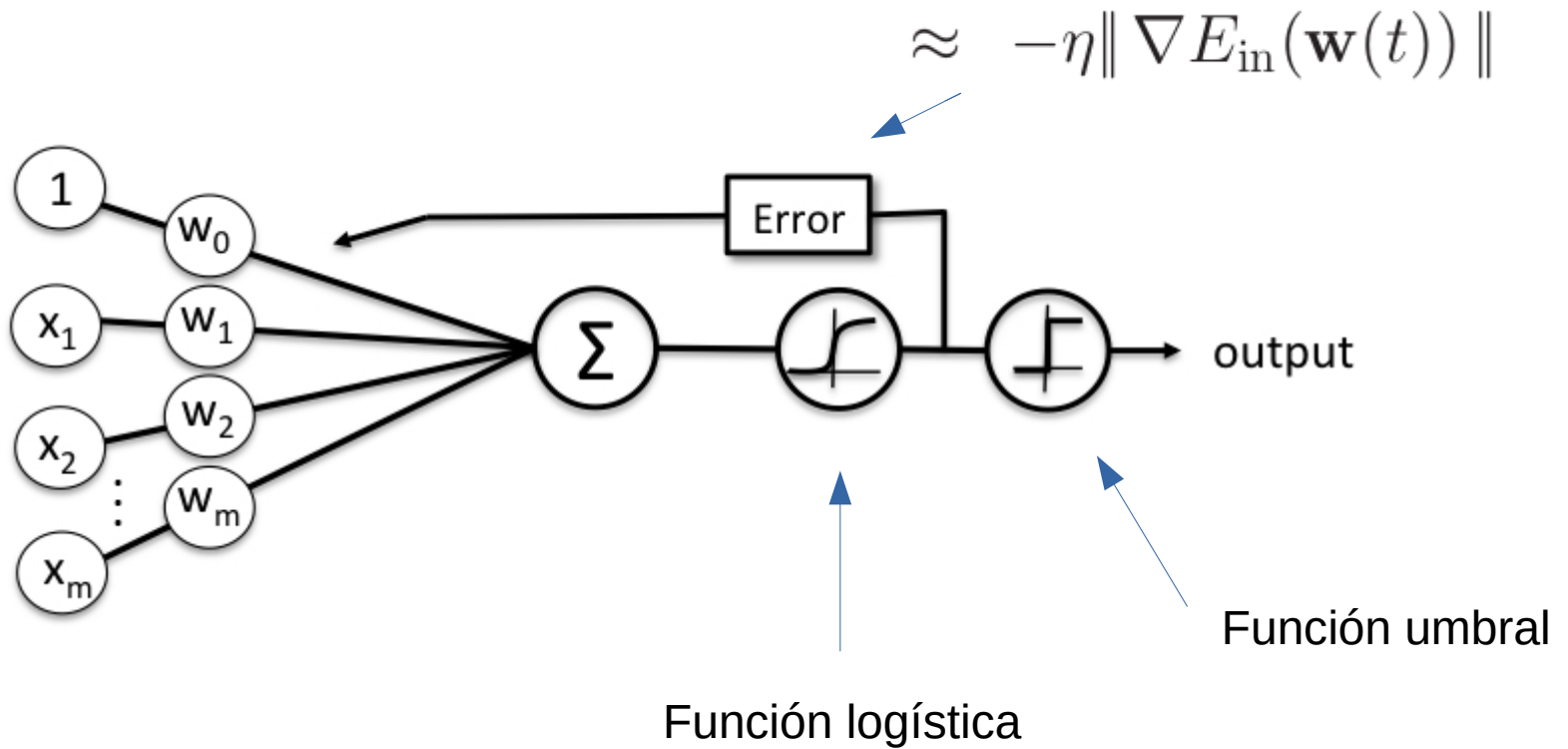
### Logistic regression algorithm:

- 1: Initialize the weights at time step  $t = 0$  to  $\mathbf{w}(0)$ .
- 2: **for**  $t = 0, 1, 2, \dots$  **do**
- 3:     Compute the gradient

$$\mathbf{g}_t = -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T(t) \mathbf{x}_n}}.$$

- 4:     Set the direction to move,  $\mathbf{v}_t = -\mathbf{g}_t$ .
- 5:     Update the weights:  $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \mathbf{v}_t$ .
- 6:     Iterate to the next step until it is time to stop.
- 7: Return the final weights  $\mathbf{w}$ .

## Regresión logística + gradiente descendente

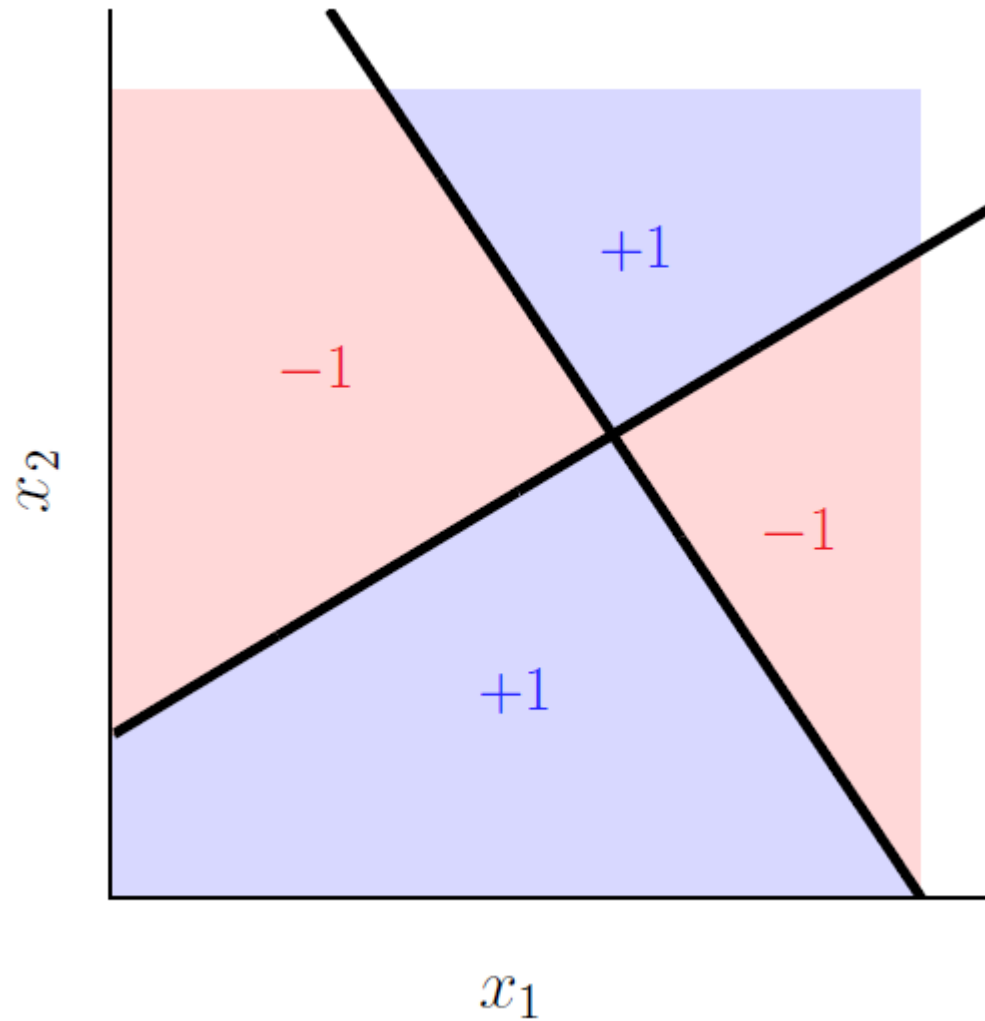




- MLP -

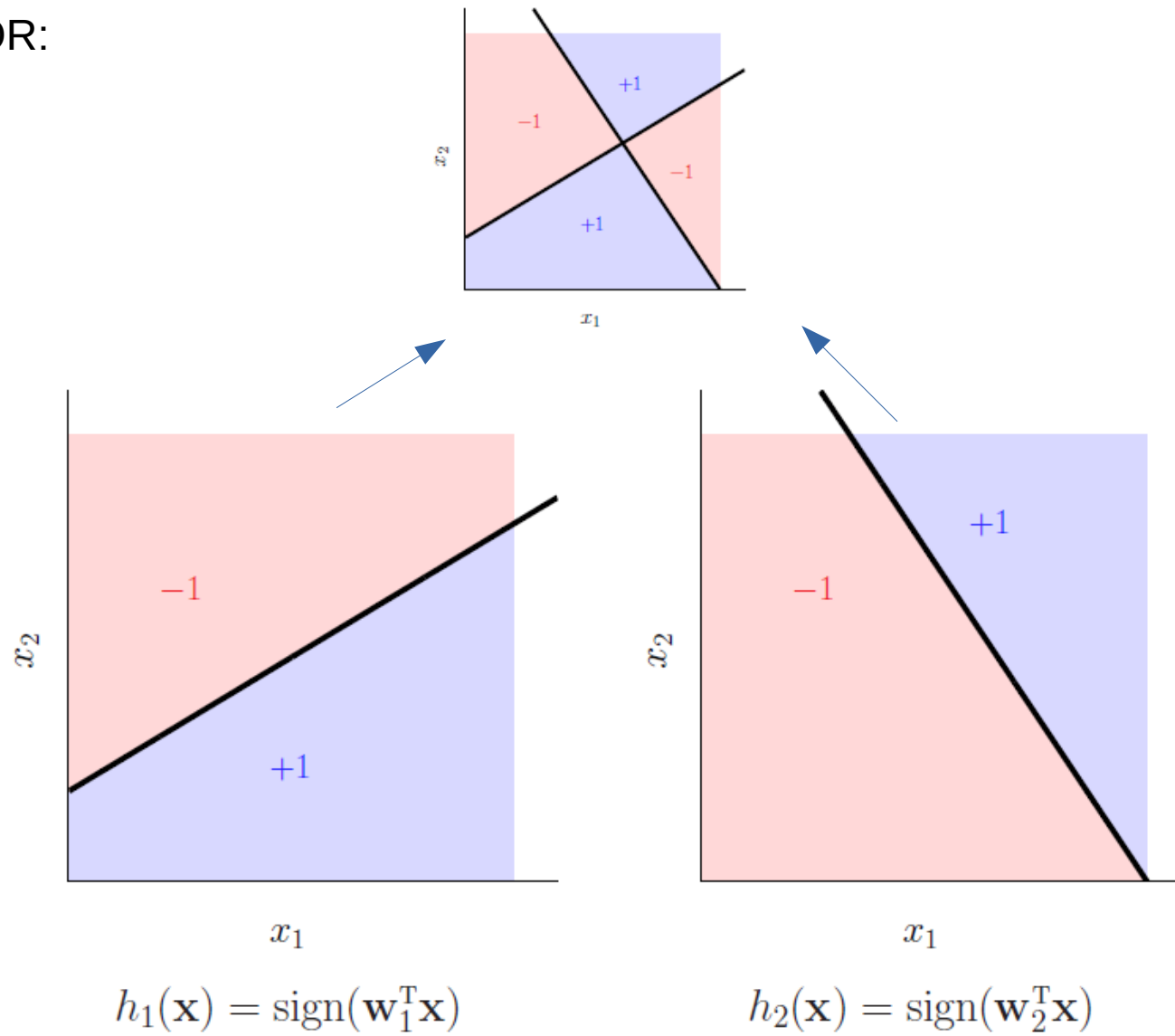
## La limitación del modelo lineal

XOR:



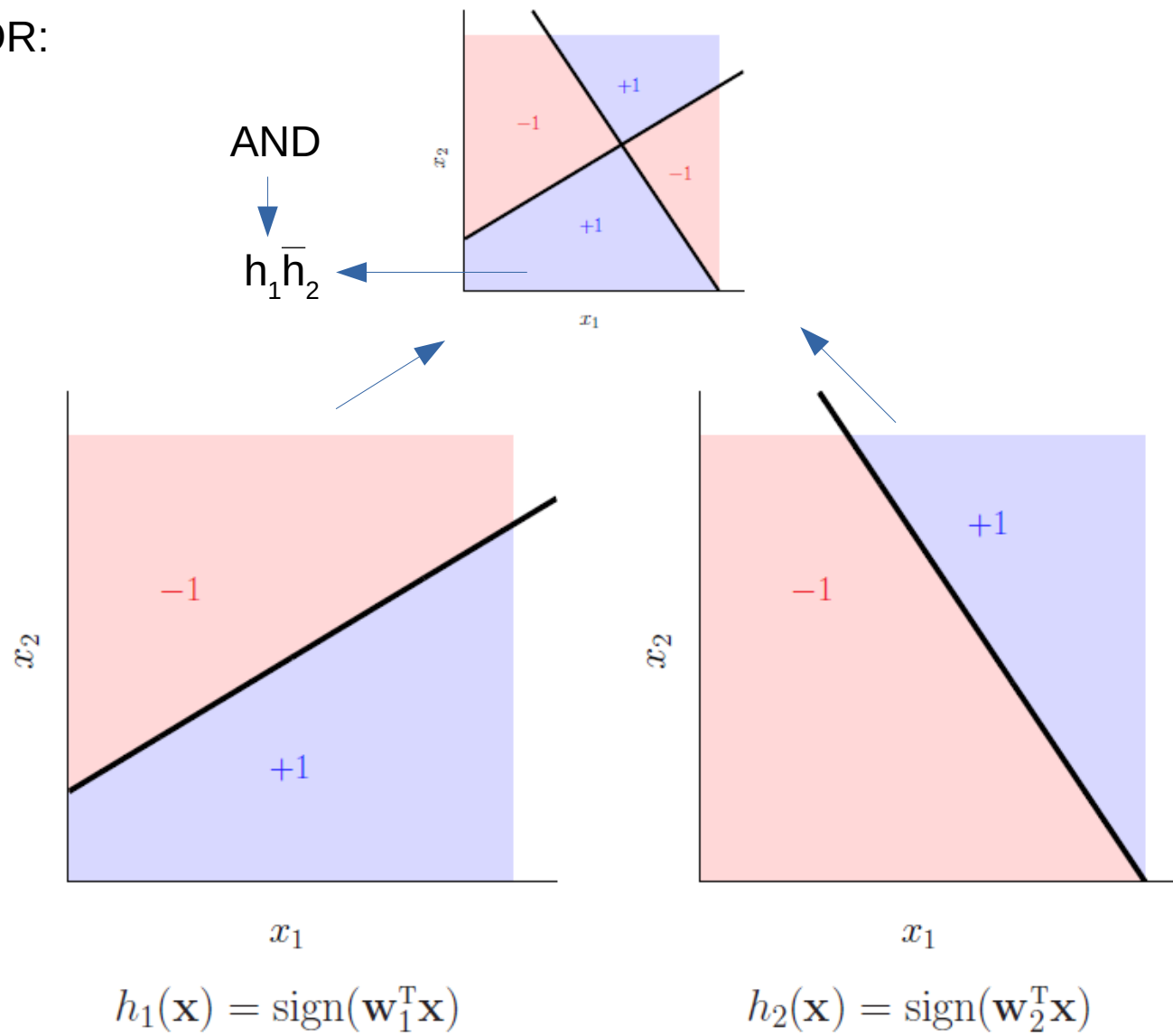
# Combinación de perceptrones

XOR:



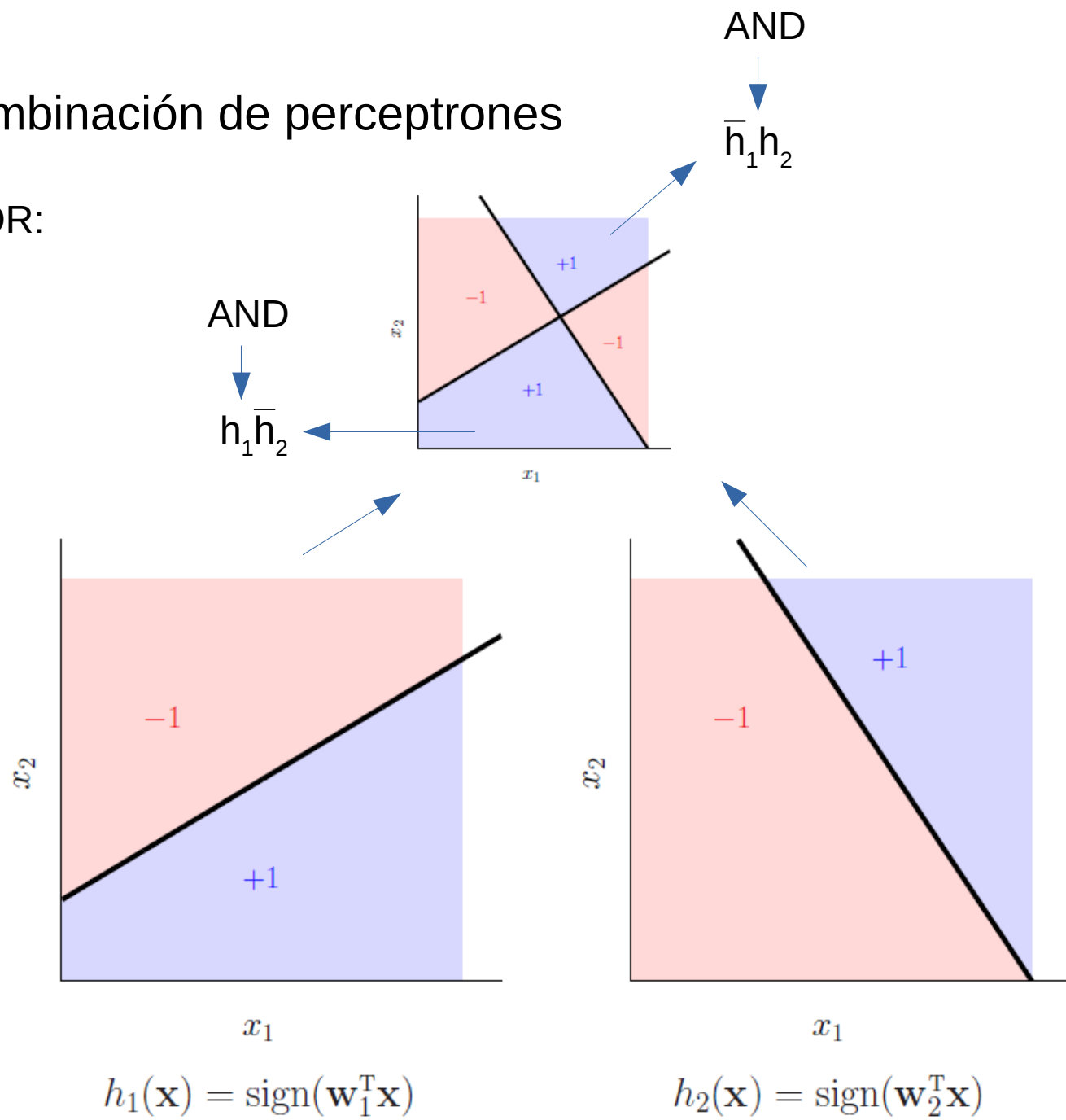
# Combinación de perceptrones

XOR:



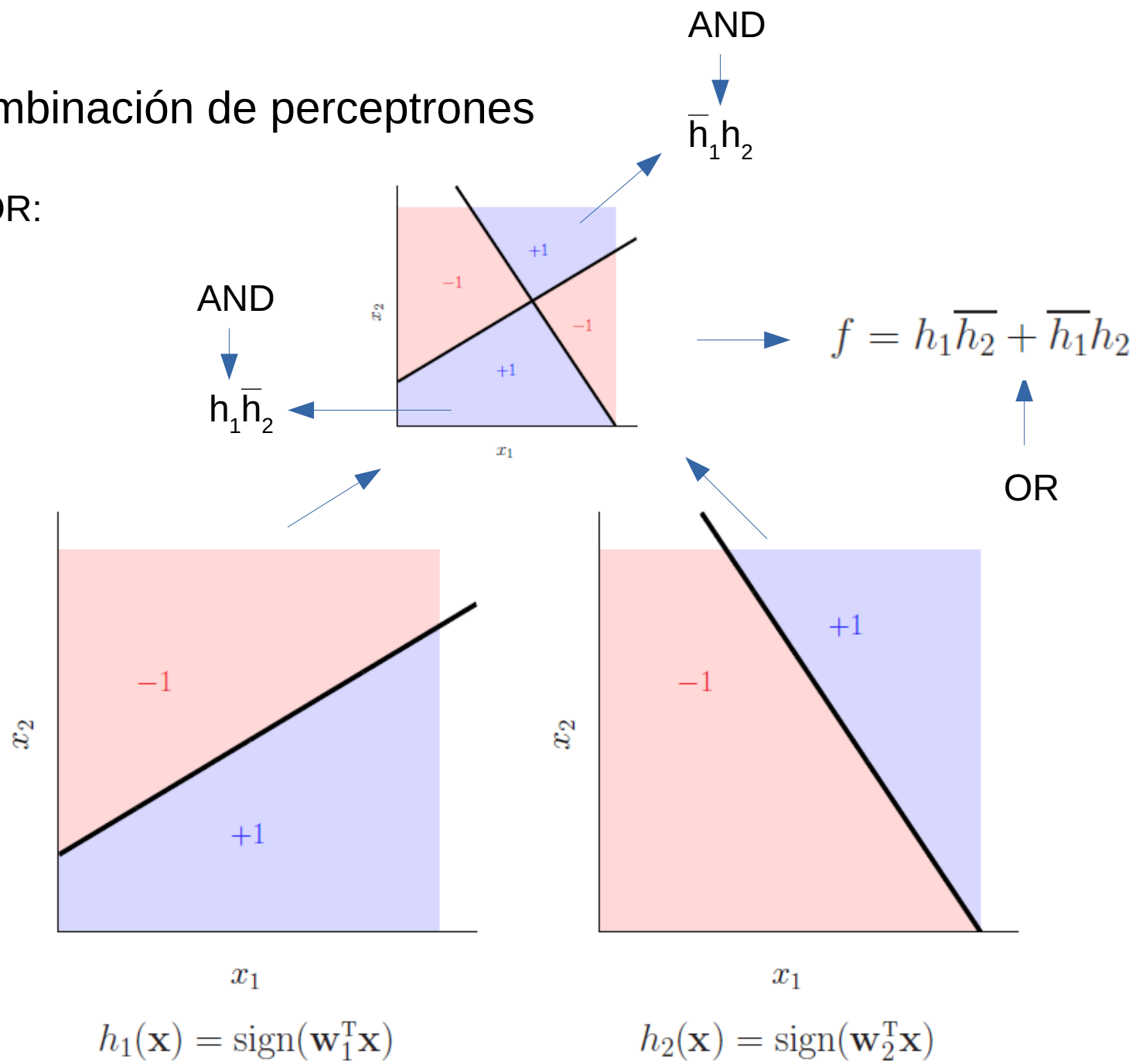
# Combinación de perceptrones

XOR:



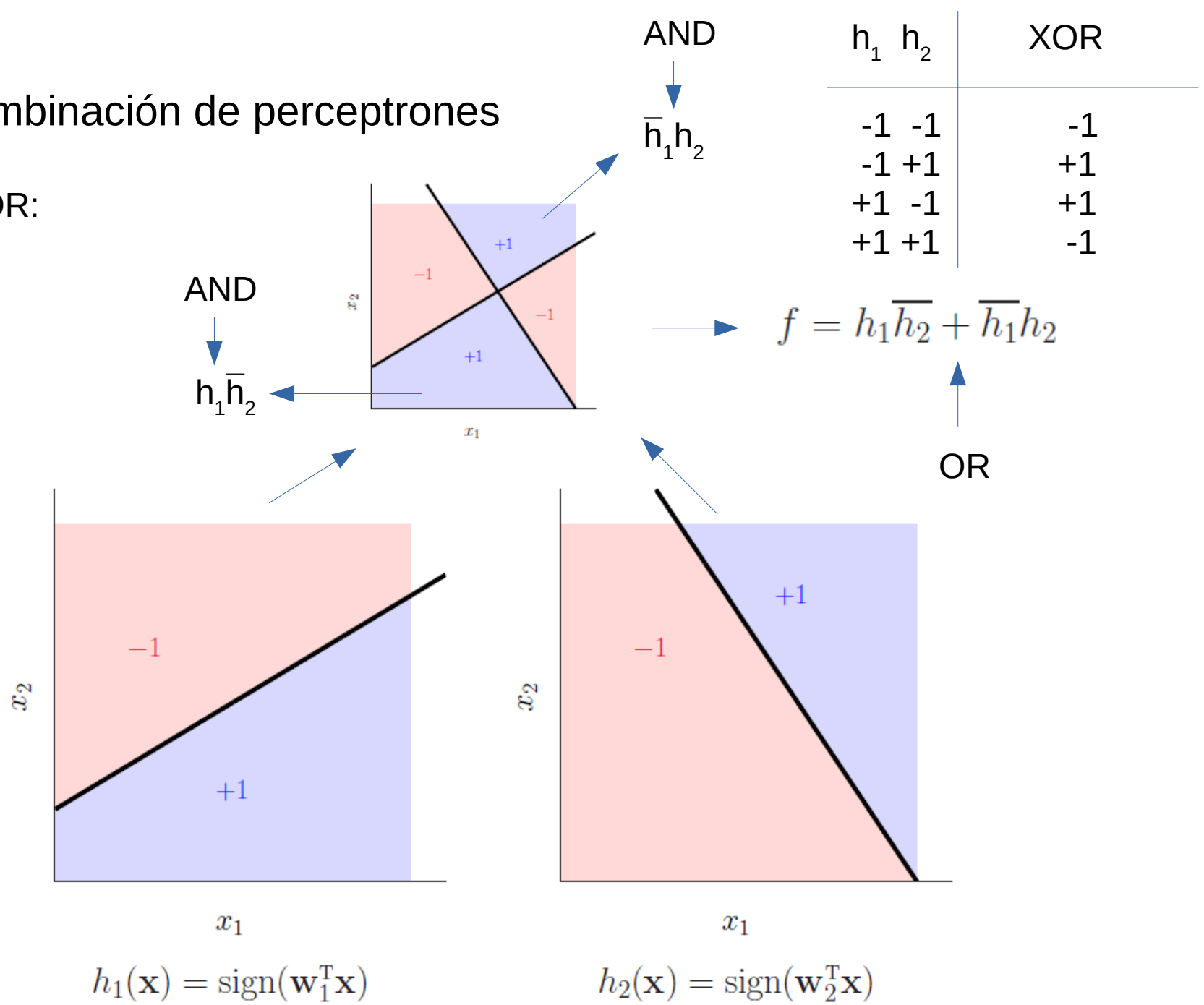
# Combinación de perceptrones

XOR:



# Combinación de perceptrones

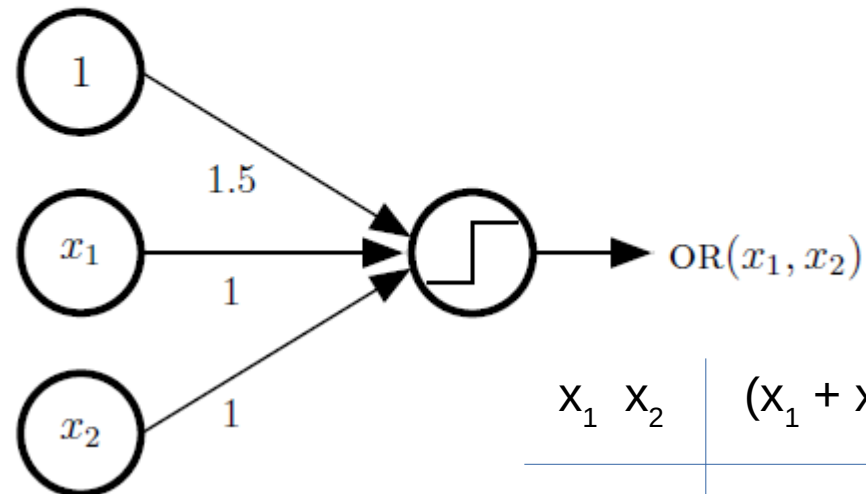
XOR:



## Combinación de perceptrones

OR:

$$\text{OR}(x_1, x_2) = \text{sign}(x_1 + x_2 + 1.5)$$



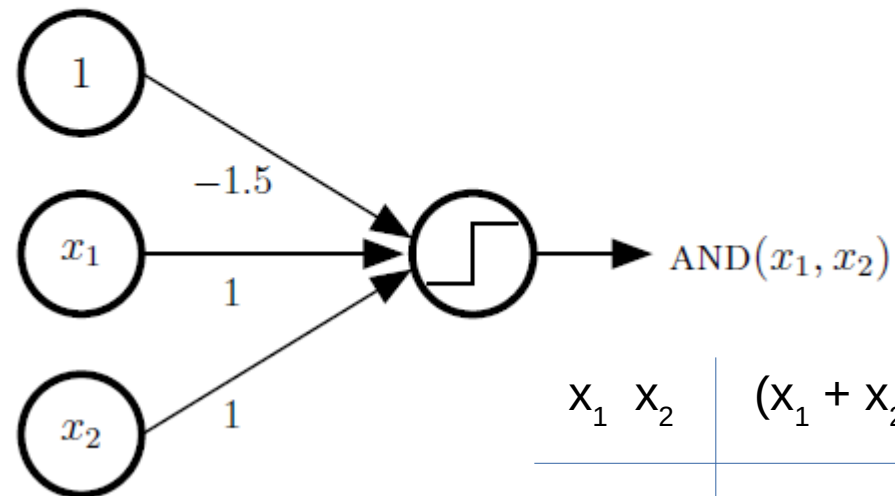
$x_1$	$x_2$	$(x_1 + x_2 + 1.5)$	signo
-1	-1	-0.5	-1
-1	+1	1.5	+1
+1	-1	1.5	+1
+1	+1	3.5	+1



## Combinación de perceptrones

AND:

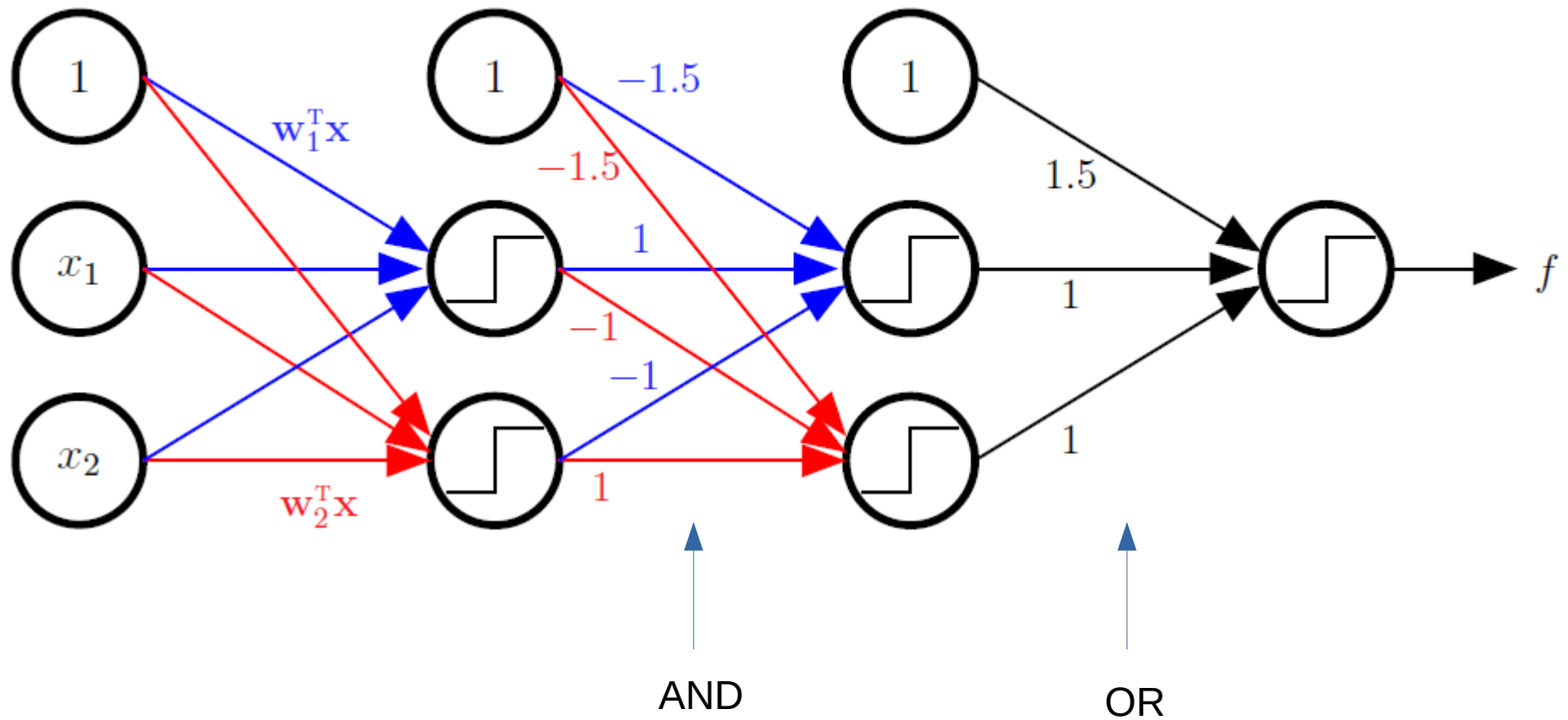
$$\text{AND}(x_1, x_2) = \text{sign}(x_1 + x_2 - 1.5)$$



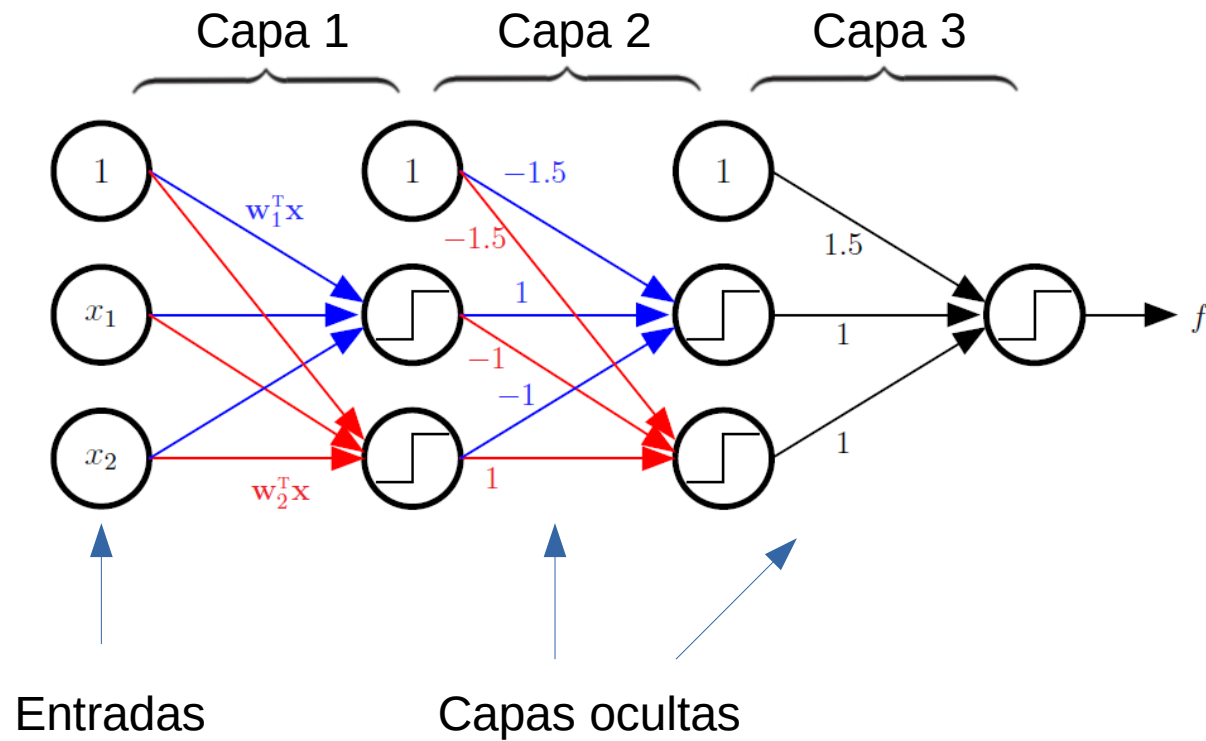
$x_1$	$x_2$	$(x_1 + x_2 - 1.5)$	signo
-1	-1	-3.5	-1
-1	+1	-1.5	-1
+1	-1	-1.5	-1
+1	+1	0.5	+1

## Combinación de perceptrones

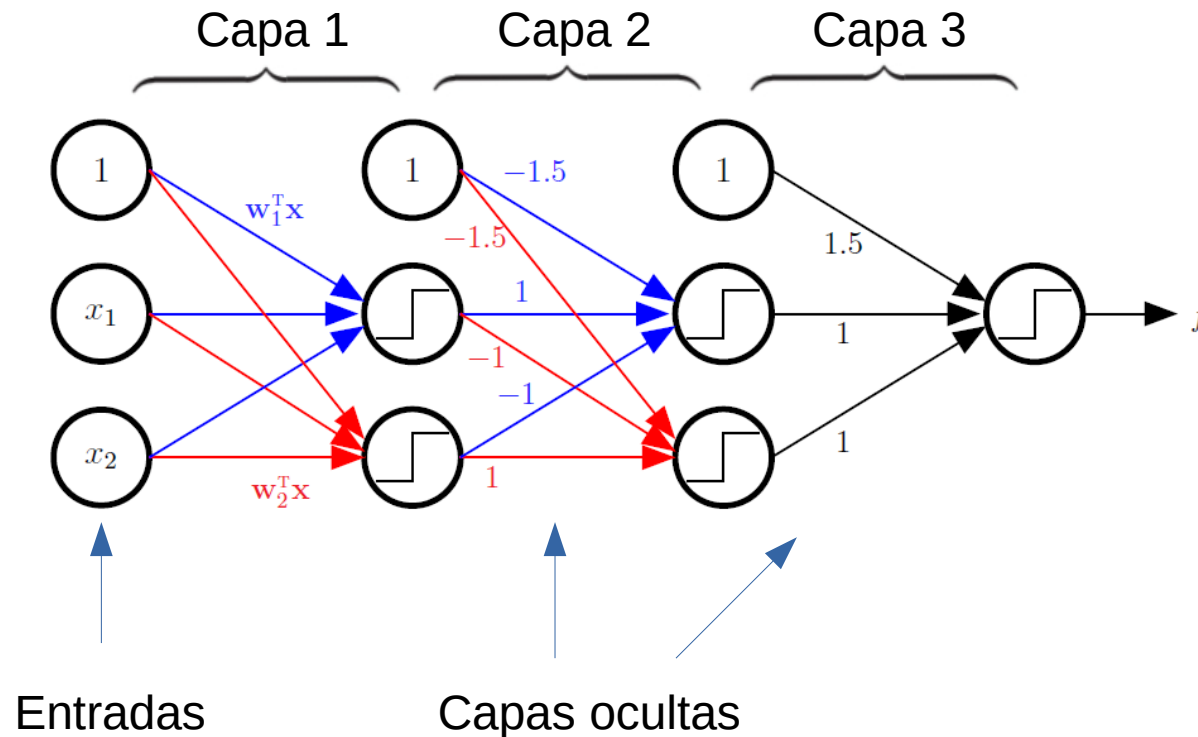
$$f = h_1 \overline{h_2} + \overline{h_1} h_2$$



# El perceptrón multicapa (Multi-Layer Perceptron)



# El perceptrón multicapa (Multi-Layer Perceptron)

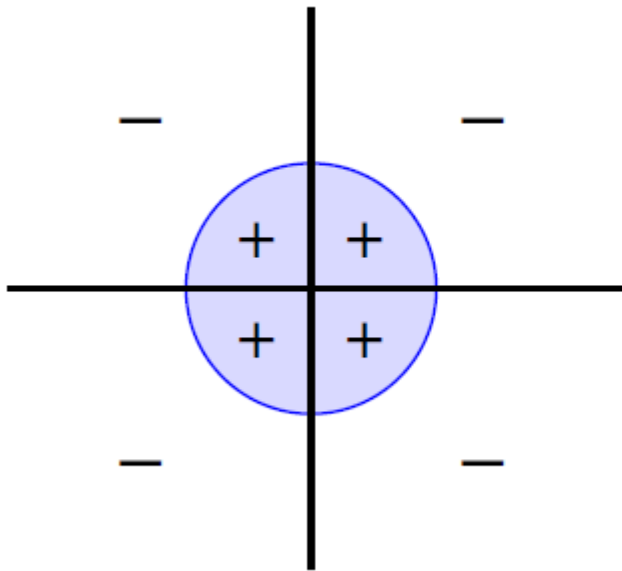


## Aproximación universal

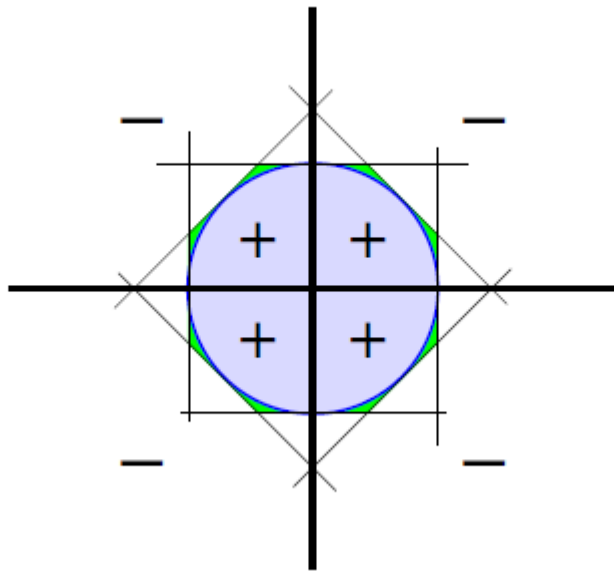
Cualquier función que se puede descomponer en separadores lineales puede ser implementada por un MLP de 3 capas

# El perceptrón multicapa (Multi-Layer Perceptron)

Un separador suave puede ser aproximado por N separadores lineales.



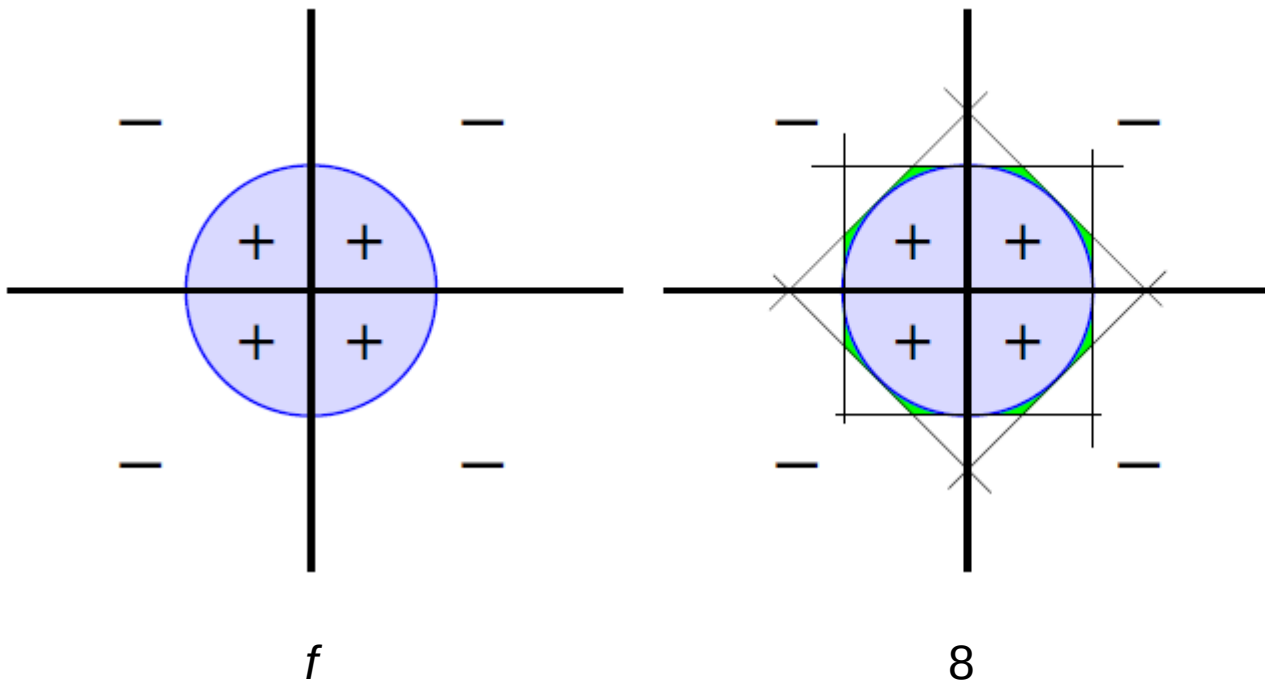
$f$



¿Cuántos perceptrones hay?

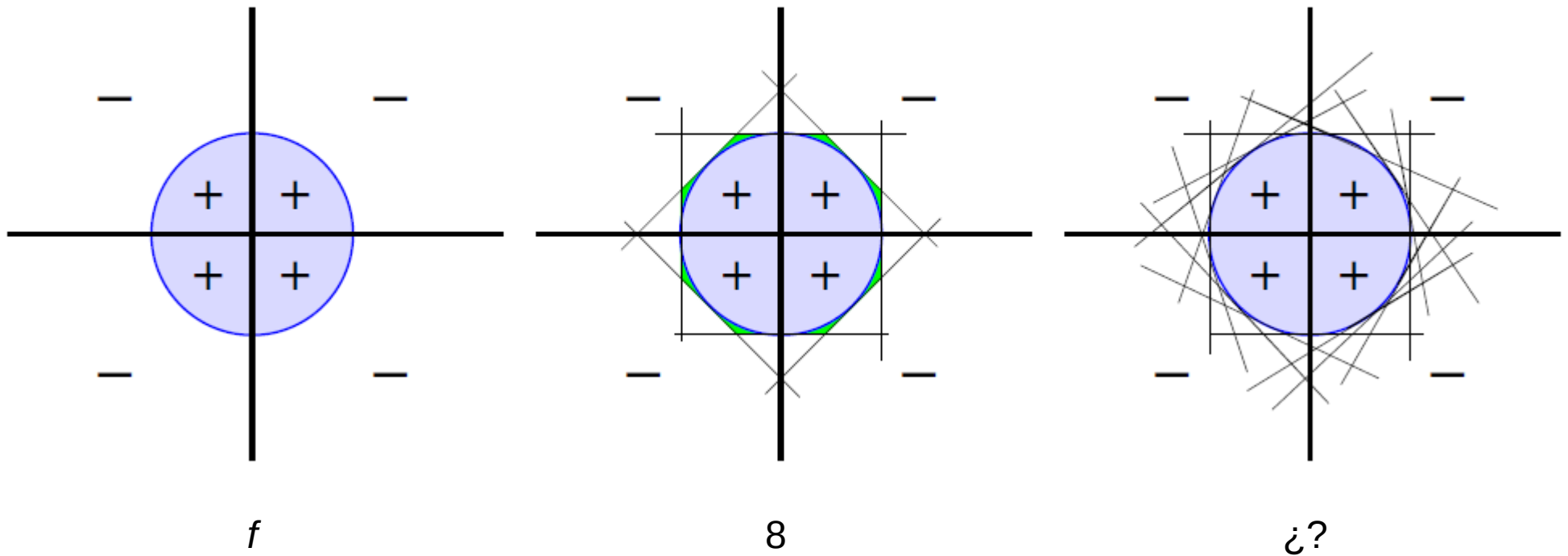
# El perceptrón multicapa (Multi-Layer Perceptron)

Un separador suave puede ser aproximado por N separadores lineales.



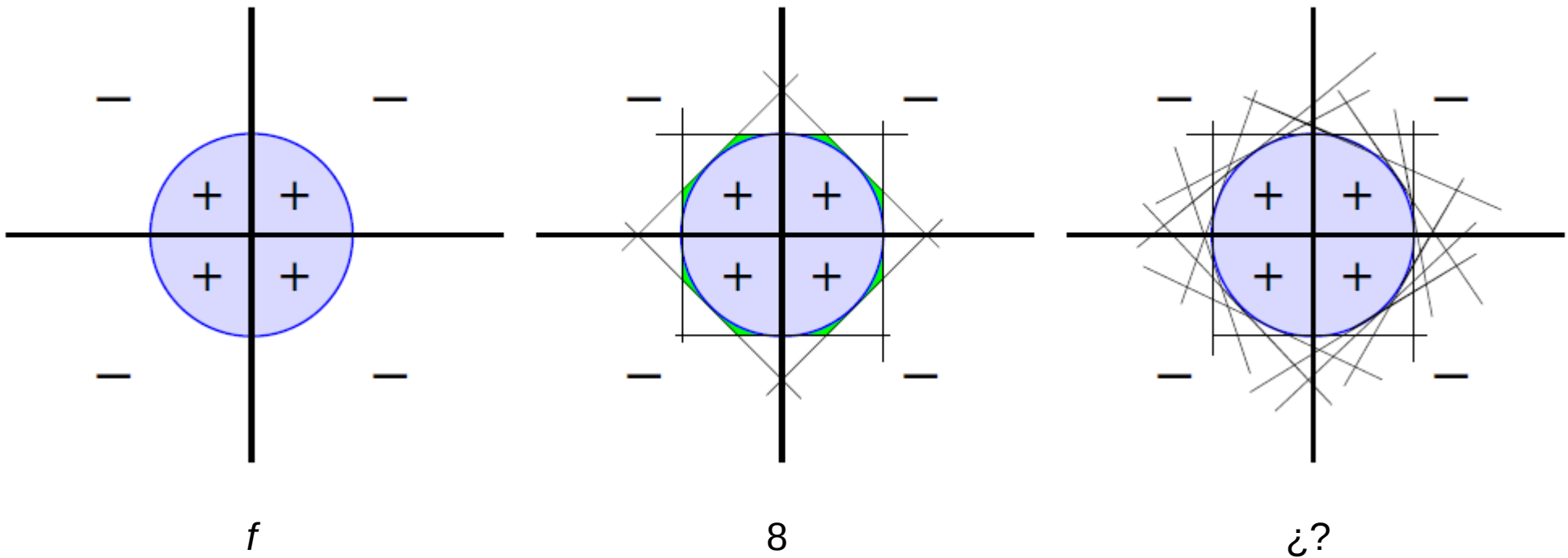
# El perceptrón multicapa (Multi-Layer Perceptron)

Un separador suave puede ser aproximado por N separadores lineales.



# El perceptrón multicapa (Multi-Layer Perceptron)

Un separador suave puede ser aproximado por N separadores lineales.



Tradeoff aproximación-generalización:

Más neuronas, mejor aproximación.

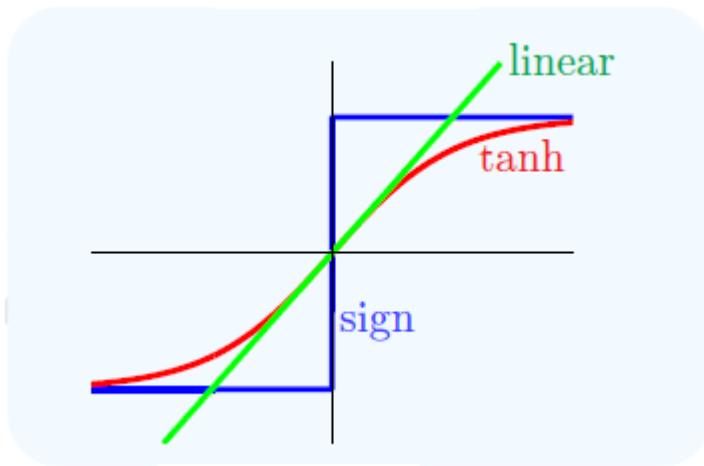
Más neuronas, peor generalización.



# El perceptrón multicapa (Multi-Layer Perceptron)

Para entrenar un MLP necesitamos reemplazar la función signo dado que no es diferenciable.

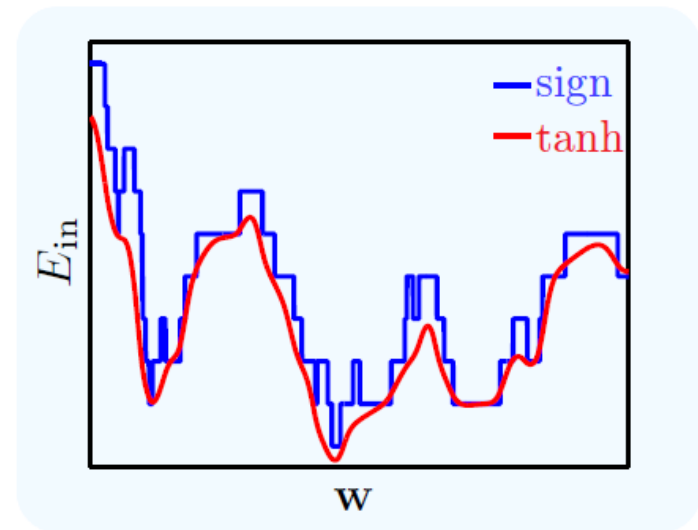
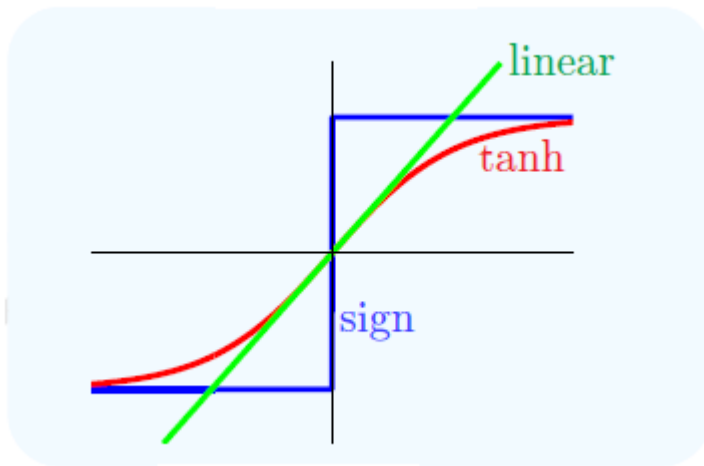
Podemos suavizar la función signo con la tangente hiperbólica, que tiene un comportamiento lineal cerca del origen y es cercana a +1 o -1 para entradas grandes.



# El perceptrón multicapa (Multi-Layer Perceptron)

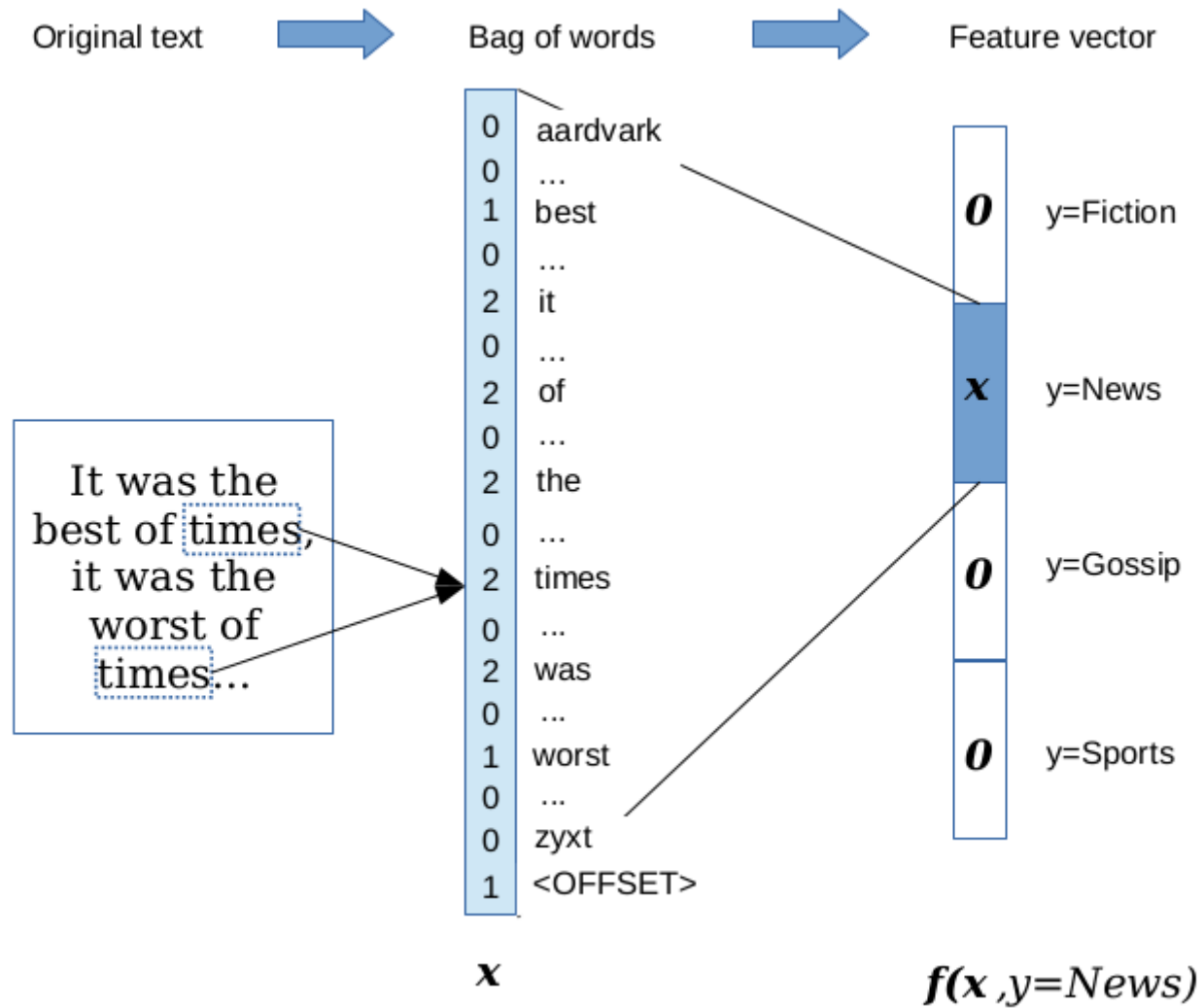
Para entrenar un MLP necesitamos reemplazar la función signo dado que no es diferenciable.

Podemos suavizar la función signo con la tangente hiperbólica, que tiene un comportamiento lineal cerca del origen y es cercana a +1 o -1 para entradas grandes.



## - CLASIFICACIÓN DE TEXTO -

# BOW



## Clasificación con word embeddings

```
inputs = Input(shape=(max_tokens, ))

embeddings_layer = Embedding(input_dim=len(tokenizer.index_word)+1,
                              output_dim=embed_len, input_length=max_tokens, trainable=True)

dense1 = Dense(128, activation="relu")
dense2 = Dense(64, activation="relu")
dense3 = Dense(len(classes), activation="softmax")

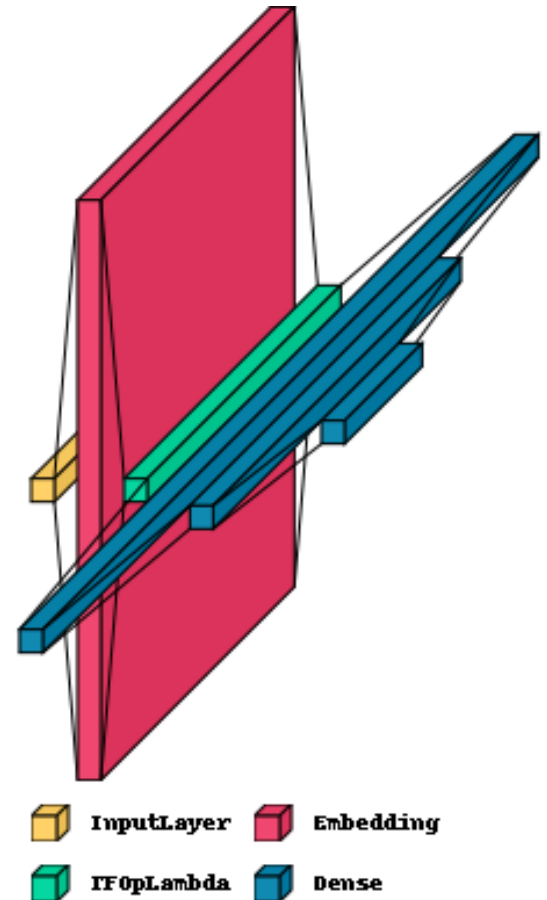
x = embeddings_layer(inputs)
x = tensorflow.reduce_sum(x, axis=1)
x = dense1(x)
x = dense2(x)
outputs = dense3(x)

model = Model(inputs=inputs, outputs=outputs)
```

| V |



forward



## Clasificación con word embeddings

