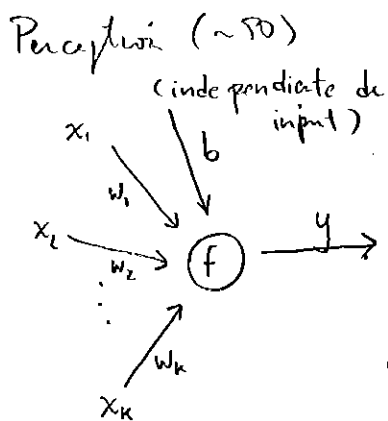
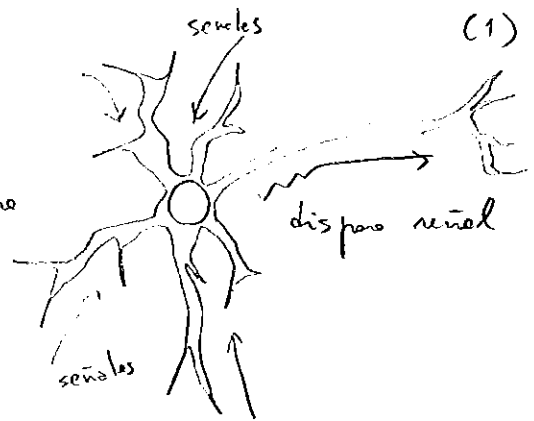


# 1) Intro a redes neuronales



motivados por neurona



$$u = x_1 w_1 + x_2 w_2 + \dots + x_k w_k + b$$

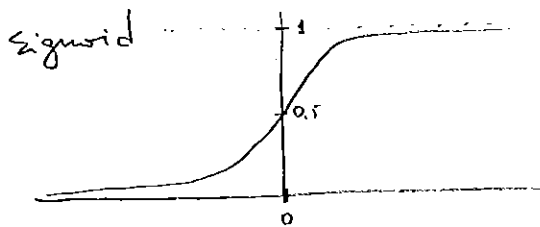
$$y = f(u) \quad \leftarrow f \text{ es misma función de activación}$$

$w_1, w_2, \dots, w_k$  : parámetros

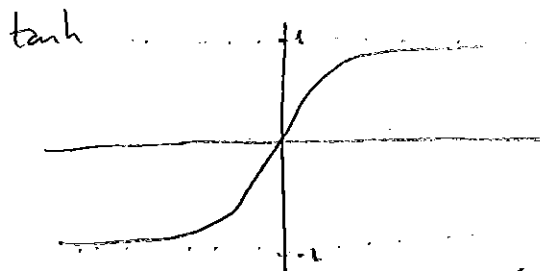
$b$  : sesgo (bias)

$x_1, x_2, \dots, x_k$  valores de entrada

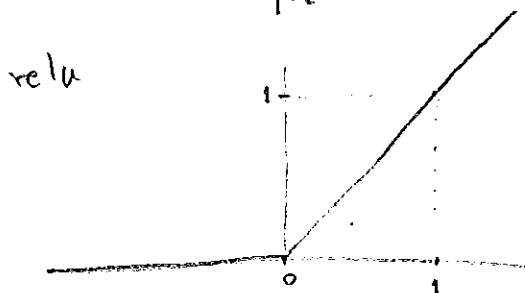
Función de activación



$$\text{sig}(u) = \frac{1}{1 + e^{-u}}$$



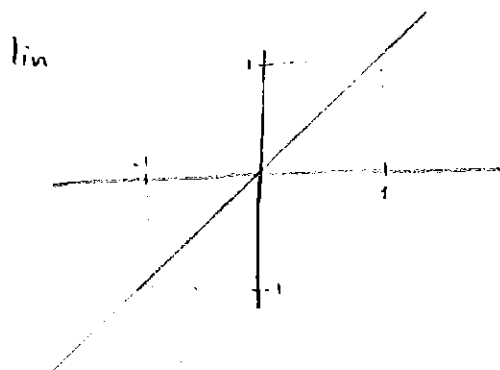
$$\text{tanh}(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}}$$



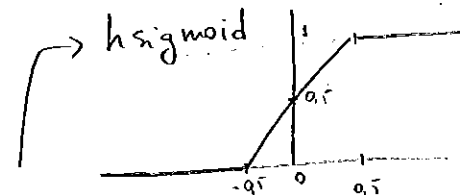
$$\text{relu}(u) = \begin{cases} u & \text{si } u > 0 \\ 0 & \text{si } u \leq 0 \end{cases}$$



Lo más importante / útil !!!

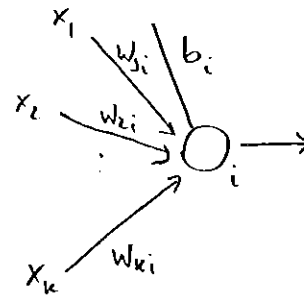
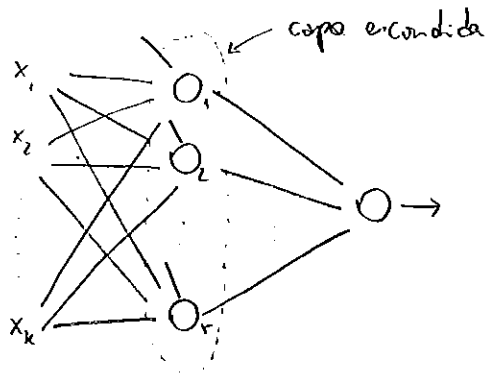


$$\text{lin}(u) = u$$



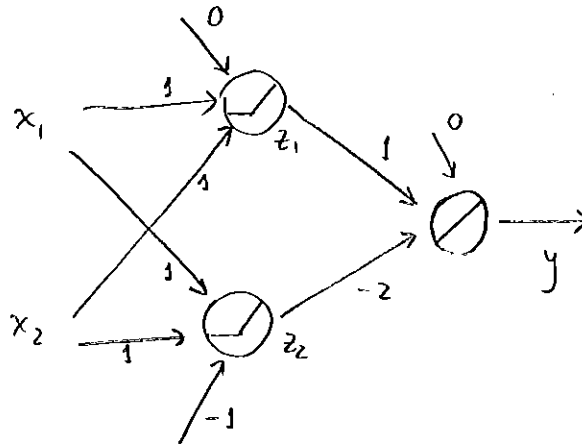
Hay otras funciones de activación  
(tendrán que usar algunas en la tarea)

redes perceptrones:



Ejemplo: XOR

$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0



$$h_1 = \text{relu}(x_1 \cdot 1 + x_2 \cdot 1 + 0) \quad y = h_1 - 2h_2$$

$$h_2 = \text{relu}(x_1 \cdot 1 + x_2 \cdot 1 - 1)$$

$$\begin{matrix} x_1 = 0 \\ x_2 = 0 \end{matrix} \Rightarrow \begin{matrix} h_1 = \text{relu}(0 + 0) = 0 \\ h_2 = \text{relu}(0 + 0 - 1) = \text{relu}(-1) = 0 \end{matrix} \quad y = 0 - 2 \cdot 0 = 0 \quad \checkmark$$

$$\begin{matrix} x_1 = 0 \\ x_2 = 1 \end{matrix} \Rightarrow \begin{matrix} h_1 = \text{relu}(0 + 1) = 1 \\ h_2 = \text{relu}(0 + 1 - 1) = 0 \end{matrix} \quad y = 1 - 2 \cdot 0 = 1 \quad \checkmark$$

$$\begin{matrix} x_1 = 1 \\ x_2 = 0 \end{matrix} \Rightarrow \begin{matrix} h_1 = \text{relu}(1 + 0) = 1 \\ h_2 = \text{relu}(1 + 0 - 1) = 0 \end{matrix} \quad y = 1 - 2 \cdot 0 = 1 \quad \checkmark$$

$$\begin{matrix} x_1 = 1 \\ x_2 = 1 \end{matrix} \Rightarrow \begin{matrix} h_1 = \text{relu}(1 + 1) = \text{relu}(2) = 2 \\ h_2 = \text{relu}(1 + 1 - 1) = \text{relu}(1) = 1 \end{matrix} \quad y = 2 - 2 \cdot 1 = 0 \quad \checkmark$$

↑ Le red de perceptrons compute de une forme  
distante a como sistema de computación,  
no "manipula números", ni no opera sobre

puede escribirse así

$$(h_1 \ h_2) = \text{relu} \left( (x_1 \ x_2) \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} + (0 \ -1) \right)$$

$$(y) = (h_1 \ h_2) \begin{pmatrix} 1 \\ -2 \end{pmatrix}$$

$$x = (x_1 \ x_2) \quad W = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \quad b = (0 \ -1)$$

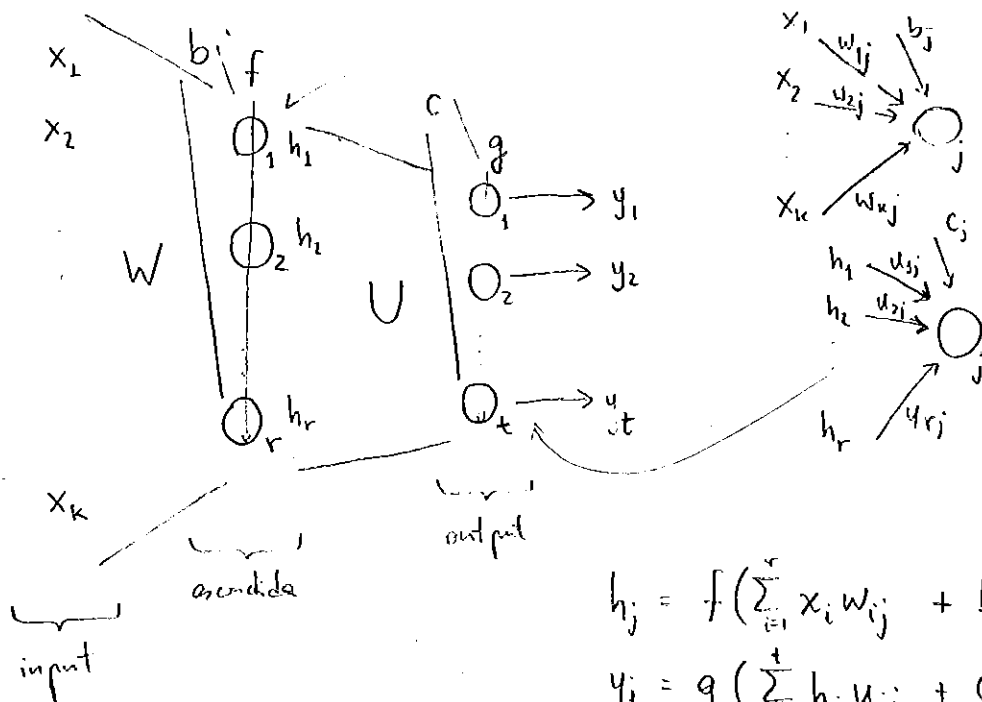
$$u = xW + b$$

$$h = \text{relu}(u)$$

$$y = hU$$

$$U = \begin{pmatrix} 1 \\ -2 \end{pmatrix}$$

Esto forma la red neuronal generalizada



$$h_j = f \left( \sum_{i=1}^r x_i w_{ij} + b_j \right)$$

$$y_i = g \left( \sum_{j=1}^t h_j u_{ij} + c_i \right)$$

$$h = f(xW + b)$$

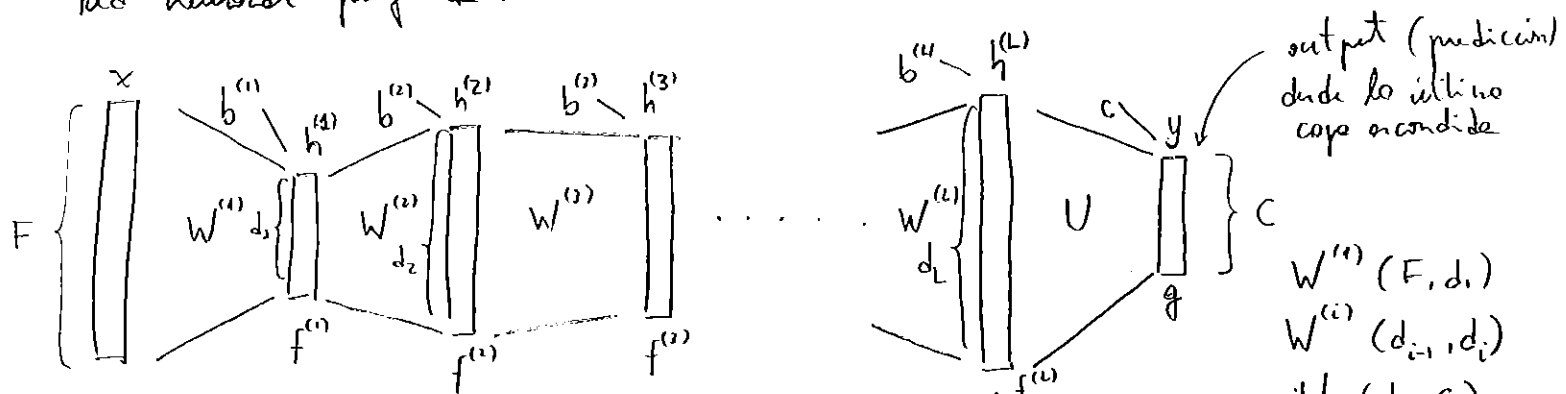
$$W_{(k,r)} \quad b_{(r)}$$

$$y = g(hU + c)$$

$$U_{(r,t)} \quad c_{(t)}$$

[ Intermedio : Universal approximation Theorem ]  $\longrightarrow$  (3b)

Red neuronal Feed Forward o Perceptrón Multicapa, o simplemente Red neuronal profunda:



$$h^{(1)} = f^{(1)}(x W^{(1)} + b^{(1)})$$

$$h^{(i)} = f^{(i)}(h^{(i-1)} W^{(i)} + b^{(i)})$$

$$y = g(h^{(L)} U + c)$$

pasado hacia adelante

$$y = \text{forward}(x) \quad (\hat{y} = \text{forward}(x))$$

Note: en general usaremos la red para computar una predicción y compararla con el valor real que yo conozco. En este caso generalmente denotaremos por  $\hat{y} = \text{forward}(x)$  para diferenciación "y" de "y" que usamos para denotar el valor "real".

¿qué usamos como función de salida? dependerá mucho de la tarea (y función de pérdida que tenemos en mente).

En general queremos probabilidades

1) clasificación binaria:

- una única neurona
- valor de salida entre 0 y 1
- Interpretación:  
 $\text{prob} \geq 0.5 \Rightarrow \text{respuesta } 1$   
 $\text{prob} < 0.5 \Rightarrow \text{respuesta } 0$

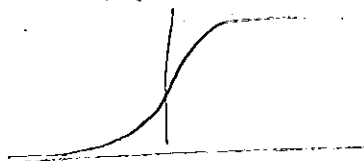
$\Rightarrow$  podemos usar sigmoid como salida

2) clasificación en varias clases

$$\hat{y} = \text{forward}(x)$$

$$P(\text{clase de } x \text{ es } 1) = \hat{y}$$

$$P(\text{clase de } x \text{ es } 0) = 1 - \hat{y}$$



(contine)

2) clasificación en varias clases:

- si tenemos  $C$  clases usamos  $C$  neuronas en la capa de salida

- necesitamos una función que tome valores

$$(u_1, u_2, u_3, \dots, u_C)$$

y me permite interpretar los como probabilidades

- lo usamos de (nuestro delante) neuronas por neurona en "softmax"

$$\text{softmax}(u_1, u_2, u_3, \dots, u_C) = (s_1, s_2, \dots, s_C)$$

$$s_i = \frac{e^{u_i}}{\sum_{j=1}^C e^{u_j}}$$

$$\text{tenemos que } \sum s_i = 1$$

$$0 \leq s_i \leq 1$$

softmax amplifica las diferencias en los valores de input

$$\text{softmax}(0, 2) = (0.12, 0.88)$$

$$\text{softmax}(0, 10) = (0.00005, 0.99995)$$

$$\text{softmax}(0, 0, 10) = (0.00005, 0.00005, 0.9999)$$

$$\text{Interpretación: } \hat{y} = \text{softmax}(u_1, u_2, \dots, u_C) = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_C)$$

$$P(x \text{ pertenece a la clase } i) = \hat{y}_i$$

según la red neuronal

En ambos casos los neuronas como nodos probabilísticos (parametrizados)

Entrenamiento: Supongamos ejemplo  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$ 

(¿qué tan bien (o mal) lo está haciendo la red?)

$$\hat{y}^{(i)} = \text{forward}(x^{(i)})$$

↑

lo que la red dice

$$y^{(i)}$$

↑

lo que esperamos

$$\text{error}(\hat{y}, y^{(i)})$$

$$\text{error}(x^{(i)}, y^{(i)})$$

↑

cuantifico "qué tan lejos" están  $y^{(i)}$  e  $\hat{y}^{(i)}$ 

lo usamos indistintamente

$$L = \frac{1}{N} \sum_{i=1}^N \text{error}(\hat{y}^{(i)}, y^{(i)}) \quad \leftarrow \text{promedio de los errores que cometen sobre todos los ejemplos.} \quad (6)$$

¿de qué depende  $L$ ? depende de la red y lo que define a la red:  $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, \dots, W^{(L)}, b^{(L)}, U, c$   
 $\theta \leftarrow \text{parámetros}$

$L(\theta)$   $L$  depende de  $\theta$

me encantaría que  $L$  me tan pequeño como se pueda

$\Rightarrow$  Problema: encontrar  $\hat{\theta}$  que minimiza  $L(\theta)$ , ¡basta!  $\min_{\theta} L(\theta)$

$L$  capas,  $d$  neuronas en capa oculta (nos chito)  $\Rightarrow O(Ld^2)$

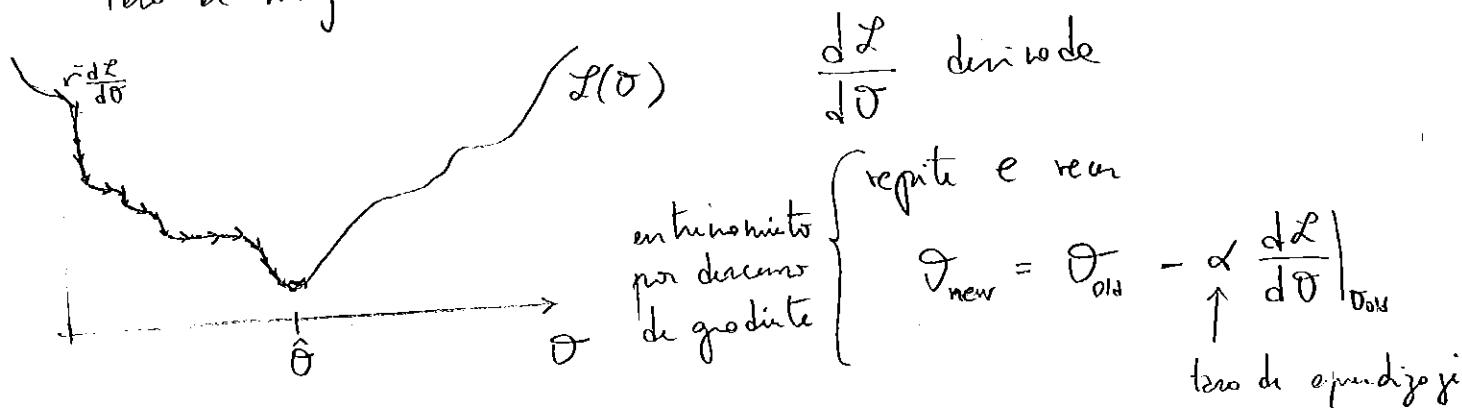
Idea: si muevo "un poquito" uno de los parámetros ¿cómo cambia la pérdida/error?

- si el error disminuye me conviene
- si el error aumenta no me conviene el cambio

$$\begin{array}{lll} L=10 & d=10 & \rightarrow 1000 \\ L=10 & d=100 & \rightarrow 100.000 \\ L=10 & d=1000 & \rightarrow 10.000.000 \end{array}$$

busquemos exhaustivos no darán resultados razonables.

Pero es muy costoso hacer un cambio y evaluar  $\Rightarrow$  lo hacemos una línea de tiempo



En general lo que tenemos que hacer es actualizar en un loop como sigue inicializar  $\hat{W}^{(1)}, \hat{b}^{(1)}, \hat{U}, \hat{c}$  de manera aleatoria, decide  $\alpha$

entrenamiento por descenso de gradiente

repite e recal (o mientras alguna condición de detención no se alcance)

para  $i$  entre 1 y  $L$

gradiente

$\hat{W}^{(i)} := \hat{W}^{(i)} - \alpha \left. \frac{\partial L}{\partial W^{(i)}} \right|_{\hat{\theta}}$

$\hat{b}^{(i)} := \hat{b}^{(i)} - \alpha \left. \frac{\partial L}{\partial b^{(i)}} \right|_{\hat{\theta}}$

$\hat{U} := \hat{U} - \alpha \left. \frac{\partial L}{\partial U} \right|_{\hat{\theta}}$

$\hat{c} := \hat{c} - \alpha \left. \frac{\partial L}{\partial c} \right|_{\hat{\theta}}$

Red  $h = f(xW + b)$

$y = g(hU + c)$  mp.  $g = \text{lin}$  (identidad)

$\Rightarrow y = f(xW + b)U + c$   $x = (x_1, x_2) \in \{0, 1\}^2$

¿que pasa si usamos  $f$  como identidad?

$\Rightarrow y = (xW + b)U + c = xWU + (bU + c)$

$\Rightarrow y = xW' + b'$

si  $f$  es la función identidad  $\Rightarrow$  la red solo puede representar una función lineal.

si  $f(u) = u \cdot m + v$  (lineal)  $\Rightarrow f(xW + b)$

$= (xW + b)m + v$

donde  $M = \begin{pmatrix} m & 0 & 0 \\ 0 & m & 0 \\ 0 & 0 & m \end{pmatrix}$   $v = (v, v, v)$

$\Rightarrow y = ((xW + b)M + v)U + c$

$= xWMU + bMU + vU + c$

$y = xW' + b'$   
 $\begin{matrix} (1) & (2) & (2,1) & (1) \end{matrix}$

$y = (x_1, x_2) \begin{pmatrix} w_1' \\ w_2' \end{pmatrix} + b'$

¿puede aprender XOR?

NO!

$(0, 0) \begin{pmatrix} w_1' \\ w_2' \end{pmatrix} + b' = 0 \Rightarrow b' = 0$

$(0, 1)$  y  $(1, 0) \Rightarrow w_1' = 1, w_2' = 1$

$(1, 1) \Rightarrow w_1' + w_2' = 0 \quad (\rightarrow \leftarrow)$

ya sabemos que con solo 1 no puede  $\Rightarrow$  no linealidad es imprescindible

¿Qué clase de funciones lineales pueden aprender?

$\varphi(x_1, x_2, \dots, x_k)$  función binaria en  $k$  inputs  $x_i \in \{0, 1\}$

$$\varphi \equiv \alpha_1 \vee \alpha_2 \vee \dots \vee \alpha_r \quad \alpha_i = l_{i1} \wedge l_{i2} \wedge \dots \wedge l_{ir_i}$$

¿puedo representar  $\alpha_i$  como un perceptrón?

$$\alpha_i(x_1, x_2, x_3, x_4, x_5) = x_1 \wedge \neg x_2 \wedge \neg x_3 \wedge \neg x_4 \wedge x_5$$

$$\begin{aligned} &\leadsto x_1 + x_5 - 2x_2 - 2x_3 - 2x_4 - 1 \\ &= \begin{cases} 1 & \text{si } x_1, x_5 = 1 \text{ y } x_2, x_3, x_4 = 0 \\ \leq 0 & \text{e.o.c.} \end{cases} \quad W_i = \begin{pmatrix} 1 \\ -2 \\ -2 \\ -2 \\ 1 \end{pmatrix} \quad b = -1 \end{aligned}$$

$$\text{si } \alpha_i(x_1, \dots, x_5) = x_1 \wedge x_2 \wedge x_3 \wedge \neg x_4 \wedge \neg x_5$$

$$\begin{aligned} &\leadsto x_1 + x_2 + x_3 - 3x_4 - 3x_5 - 2 \\ &= \begin{cases} 1 & \text{si } x_1, x_2, x_3 = 1 \text{ y } x_4, x_5 = 0 \\ \leq 0 & \text{e.o.c.} \end{cases} \quad W_i = \begin{pmatrix} 1 \\ 1 \\ 1 \\ -3 \\ -3 \end{pmatrix} \quad b = -2 \end{aligned}$$

$$\alpha_i = \bigwedge_{i \in P} x_i \wedge \bigwedge_{i \in N} \neg x_i \quad \leadsto \quad \sum_{i \in P} x_i - |P| \sum_{i \in N} x_i - |P| + 1$$

$\Rightarrow$  a cada  $\alpha_i$  le puedo asignar un perceptrón con activación relé  
 cope de relé puede ser simplemente un  $\text{hsigmoid}(h_1 + h_2 + \dots + h_r - 0.5)$   
 (necesitamos que XOR necesite control,  
 exponencial de neuronas)

Por otro:

Teo (UAT): no  $f$  una función continua de  $[0, 1]^k \rightarrow [0, 1]$ .

Para todo  $\varepsilon > 0$  existen  $W, b, U$  tal la función

$$F(x) = \text{sig}(xW + b) \cup \quad \text{cumple con}$$

$$\|f - F\| < \varepsilon \quad (|f(x) - F(x)| < \varepsilon)$$

[En otras palabras, las funciones computadas por redes neuronales  
 son "densas" en el espacio de las funciones continuas  $[0, 1]^k \rightarrow [0, 1]$ ]