

# Modelos de Deep Learning

## Multi-Layer Perceptron

Universidad ORT Uruguay

1 de Setiembre, 2025

# Modelo de una neurona

- Modelo paramétrico:

$$\hat{y} = f(\mathbf{x}; \boldsymbol{\theta})$$

donde la función  $f$  está parametrizada por  $\boldsymbol{\theta}$

- Regresión lineal - Modelo paramétrico básico de regresión

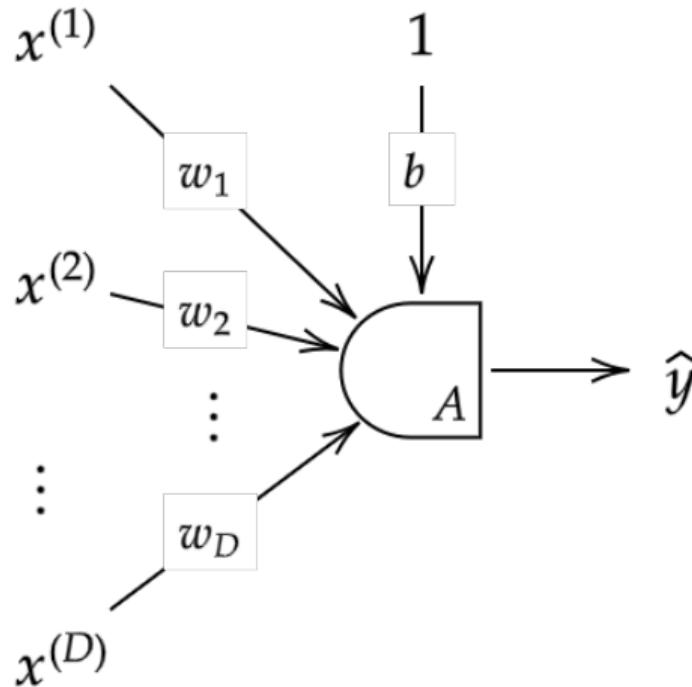
$$\hat{y} = w_1x_1 + w_2x_2 + \cdots + w_Dx_D + b = \mathbf{x}^\top \mathbf{w} + b$$

- Regresión logística - Modelo paramétrico básico de clasificación

$$\hat{y} = \text{Sigmoid}(w_1x_1 + w_2x_2 + \cdots + w_Dx_D + b) = \text{Sigmoid}(\mathbf{x}^\top \mathbf{w} + b)$$

- En ambos modelos los parámetros son  $\boldsymbol{\theta} = (\mathbf{w}, b)$

# Perceptrón: modelo de una neurona



Usando la notación matricial:

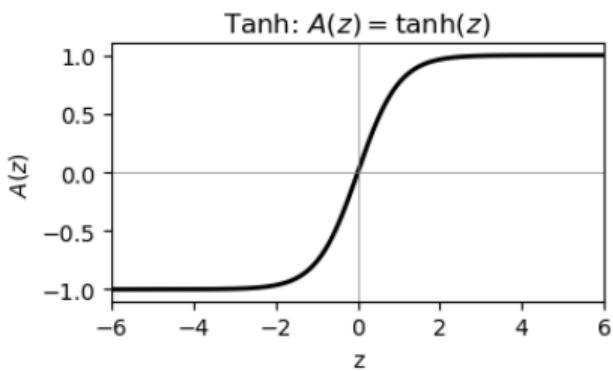
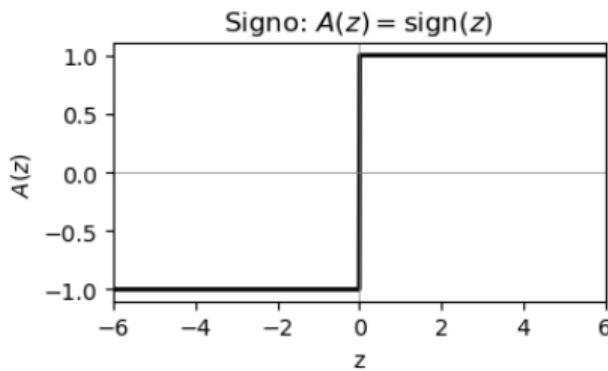
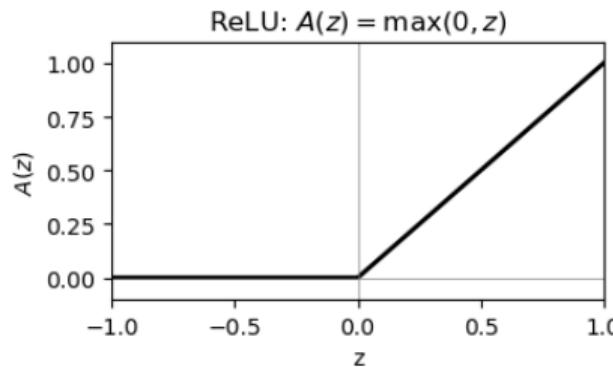
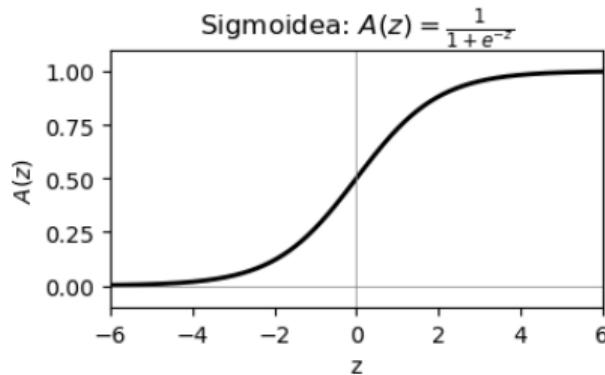
$$\hat{y} = A(\mathbf{x}^\top \mathbf{w} + b)$$

en donde

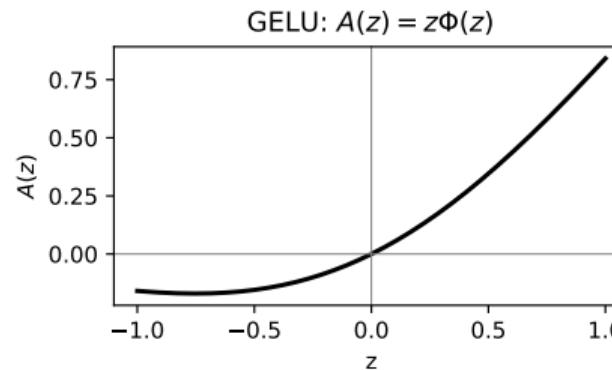
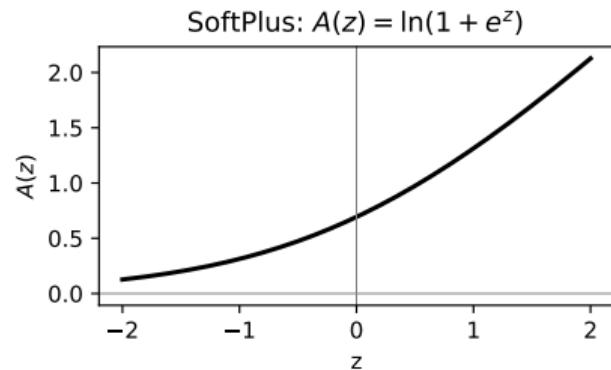
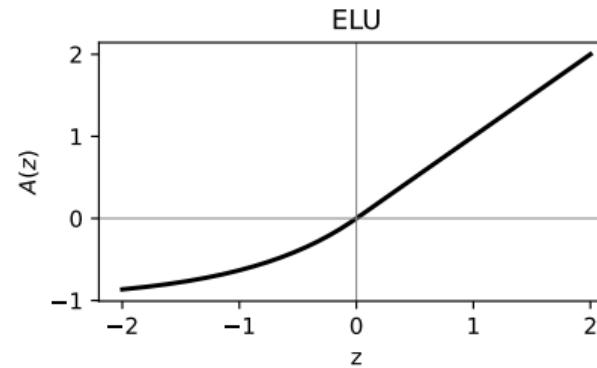
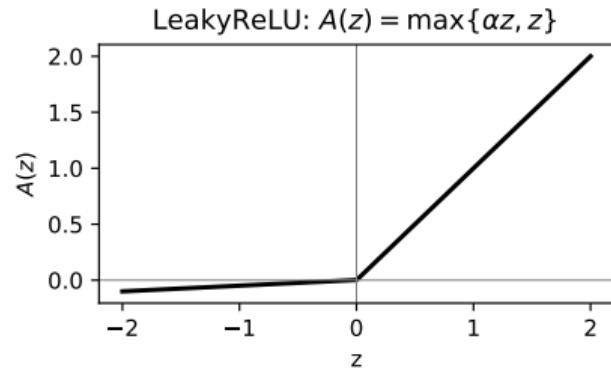
$$\mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_D \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix}$$

$A : \mathbb{R} \rightarrow \mathbb{R}$  activación

# Funciones de activación



# Funciones de activación



# Componente básica: Densa/Lineal

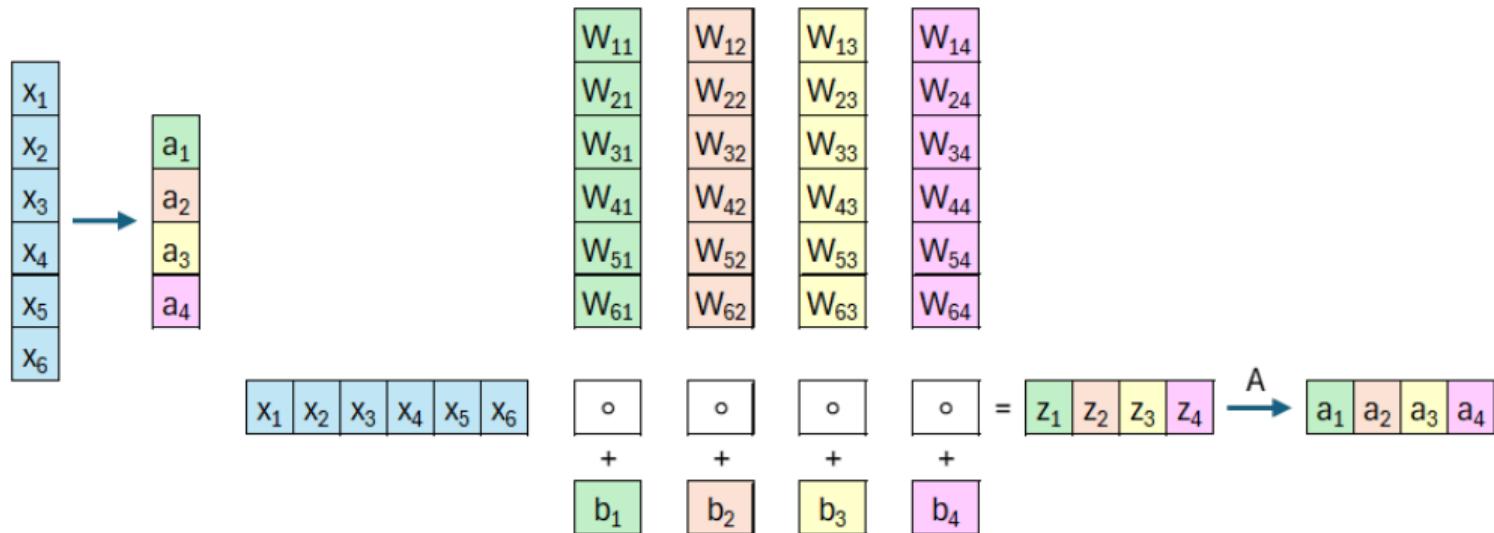
La **componente básica** de un MLP es una capa (layer) **densa/lineal**:

$$\mathbf{a} = D(\mathbf{x}; \mathbf{W}, \mathbf{b}) = A(\mathbf{x}^\top \mathbf{W} + \mathbf{b}^\top) \text{ con } \mathbf{x} \sim (I), \mathbf{a} \sim (O), \mathbf{W} \sim (I, O), \mathbf{b} \sim (O)$$

- la operación  $\mathbf{z} = \mathbf{x}^\top \mathbf{W} + \mathbf{b}^\top$  es la **unidad lineal**
- $A$  es la función de **activación** con  $\mathbf{a} = A(\mathbf{z})$  point-wise

$$\begin{aligned}\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_I \end{bmatrix} &\mapsto \mathbf{z} = [x_1 \ \cdots \ x_I] \begin{bmatrix} W_{11} & \cdots & W_{1O} \\ \vdots & \ddots & \vdots \\ W_{I1} & \cdots & W_{IO} \end{bmatrix} + [b_1 \ \cdots \ b_O] \\ &\mapsto \mathbf{a} = [a_1 \ \cdots \ a_O]\end{aligned}$$

# Ejemplo Capa Densa



La cantidad de parámetros de  $D$  es  $I \times O + O = (I + 1) \times O$ .

# Operando en batch

Usando un batch de  $N$  inputs:

- $\mathbf{X} \sim (N, I)$ ,  $\mathbf{a} \sim (N, O)$ ,  $\mathbf{W} \sim (I, O)$ ,  $\mathbf{b} \sim (O)$
- Aquí cada **fila** de  $\mathbf{X}$  es un input de shape  $(I)$ .
- A  $\mathbf{b}$  se aplica **broadcasting** para tener shape  $(N, O)$
- La **salida** de la capa es

$$\mathbf{a} = D(\mathbf{X}; \mathbf{W}, \mathbf{b}) = A(\mathbf{X}\mathbf{W} + \mathbf{b}) \sim (N, O)$$

# Ejemplo Capa Densa Operación en Batch

$W_{11}$	$W_{12}$	$W_{13}$	$W_{14}$
$W_{21}$	$W_{22}$	$W_{23}$	$W_{24}$
$W_{31}$	$W_{32}$	$W_{33}$	$W_{34}$
$W_{41}$	$W_{42}$	$W_{43}$	$W_{44}$
$W_{51}$	$W_{52}$	$W_{53}$	$W_{54}$
$W_{61}$	$W_{62}$	$W_{63}$	$W_{64}$

$$\begin{matrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} & x_4^{(1)} & x_5^{(1)} & x_6^{(1)} \end{matrix} \quad \begin{matrix} \circ \\ + \\ b_1 \end{matrix} \quad \begin{matrix} \circ \\ + \\ b_2 \end{matrix} \quad \begin{matrix} \circ \\ + \\ b_3 \end{matrix} \quad \begin{matrix} \circ \\ + \\ b_4 \end{matrix} = \begin{matrix} z_1^{(1)} & z_2^{(1)} & z_3^{(1)} & z_4^{(1)} \end{matrix} \xrightarrow{A} \begin{matrix} a_1^{(1)} & a_2^{(1)} & a_3^{(1)} & a_4^{(1)} \end{matrix}$$

$$\begin{matrix} x_1^{(2)} & x_2^{(2)} & x_3^{(2)} & x_4^{(2)} & x_5^{(2)} & x_6^{(2)} \end{matrix} \quad \begin{matrix} \circ \\ + \\ b_1 \end{matrix} \quad \begin{matrix} \circ \\ + \\ b_2 \end{matrix} \quad \begin{matrix} \circ \\ + \\ b_3 \end{matrix} \quad \begin{matrix} \circ \\ + \\ b_4 \end{matrix} = \begin{matrix} z_1^{(2)} & z_2^{(2)} & z_3^{(2)} & z_4^{(2)} \end{matrix} \xrightarrow{A} \begin{matrix} a_1^{(2)} & a_2^{(2)} & a_3^{(2)} & a_4^{(2)} \end{matrix}$$

# Multi-layer Perceptron

- Un **MLP** de  $d > 1$  *capas* se obtiene **componiendo**  $d$  *capas densas*:

$$\hat{\mathbf{y}} = f(\mathbf{x}; \mathbf{W}_1, \dots, \mathbf{W}_d, \mathbf{b}_1, \dots, \mathbf{b}_d) = D_d(D_{d-1}(\dots, \mathbf{W}_{d-1}, \mathbf{b}_{d-1}); \mathbf{W}_d, \mathbf{b}_d)$$

- La dimensión de salida de  $O_i$  tiene que ser igual a la de entrada de  $I_{i+1}$
- La cantidad de **parámetros** de  $f(\mathbf{x})$  es:

$$\text{params}(f) = \sum_{i=1}^d \text{params}(D_i) = \sum_{i=1}^d (I_i + 1)O_i$$

con  $\mathbf{W}_i \sim (I_i, O_i)$ ,  $\mathbf{b}_i \sim (O_i)$ .