

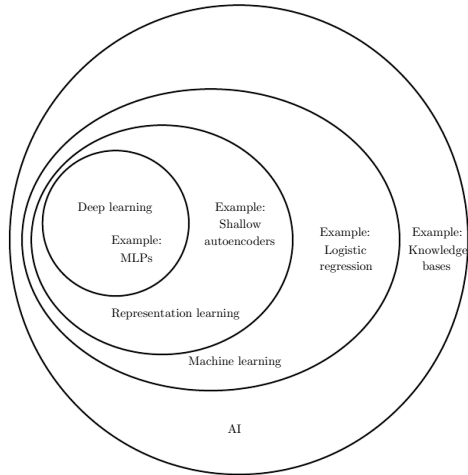
Machine Learning para Inteligencia Artificial

Redes Neuronales

Universidad ORT Uruguay

4 de Junio, 2025

Inteligencia Artificial (AI) - Machine learning (ML)



Representación

- Datos
- Hipótesis

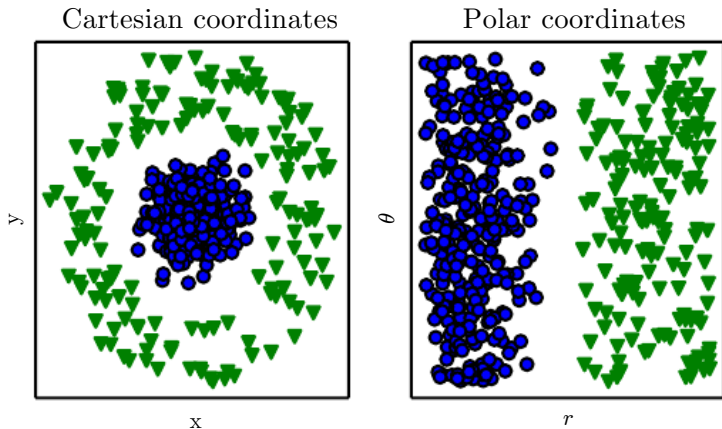
Memorización

Inferencia/generalización

I. Goodfellow, Y. Bengio, A. Courville. Deep Learning. MIT Press, 2016.

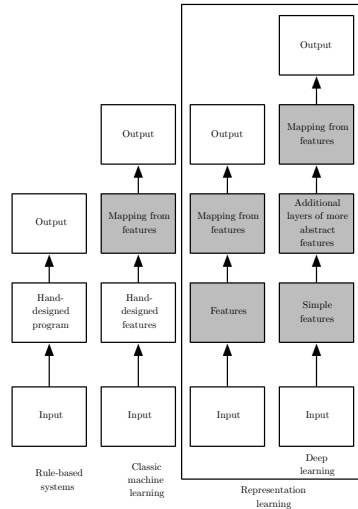
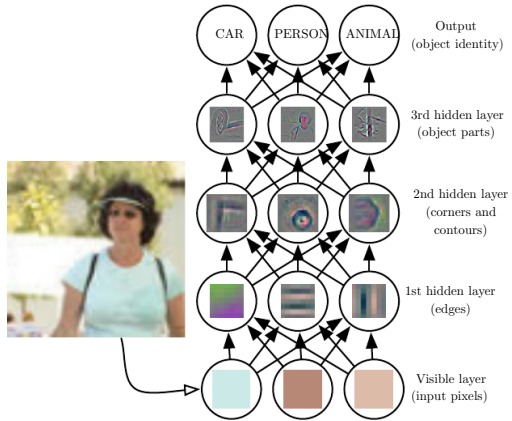
Representación

Punto \longrightarrow sistemas de coordenadas



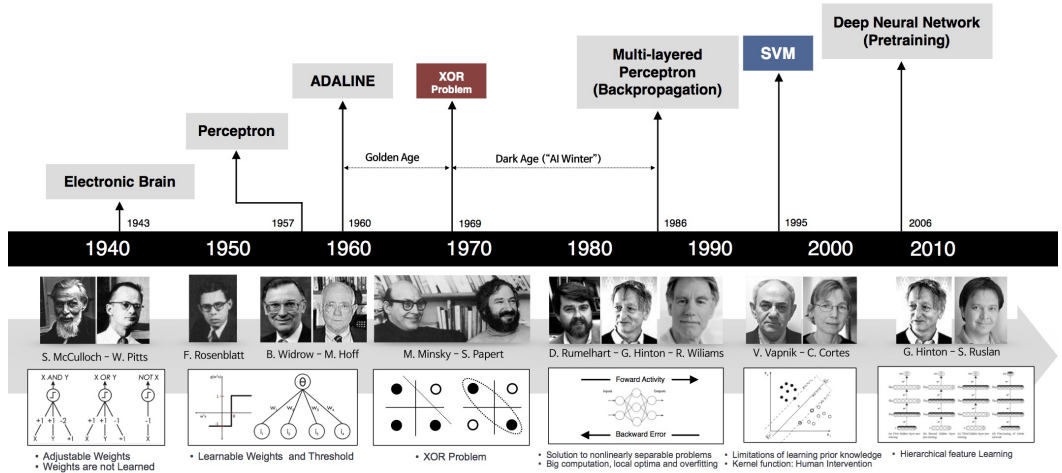
I. Goodfellow, Y. Bengio, A. Courville. Deep Learning. MIT Press, 2016.

Representación composicional

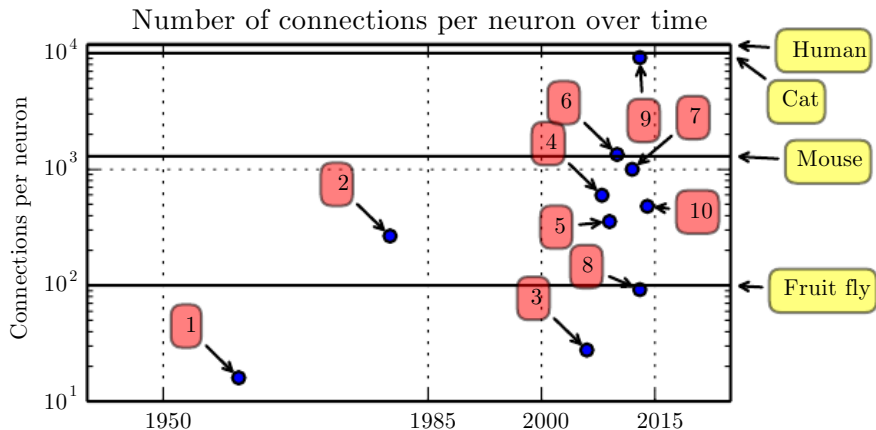


I. Goodfellow, Y. Bengio, A. Courville. Deep Learning. MIT Press, 2016.

Deep Learning: origen

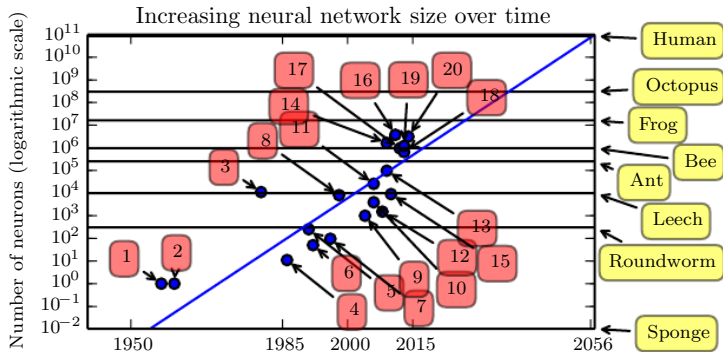


Artificial Neural Networks - Cantidad de connexiones



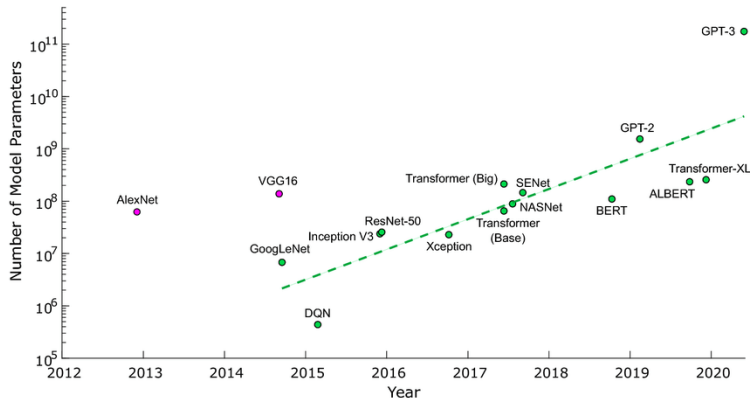
I. Goodfellow, Y. Bengio, A. Courville. Deep Learning. MIT Press, 2016.

Artificial Neural Networks - Cantidad de neuronas



I. Goodfellow, Y. Bengio, A. Courville. Deep Learning. MIT Press, 2016.

Artificial Neural Networks - Cantidad de parámetros



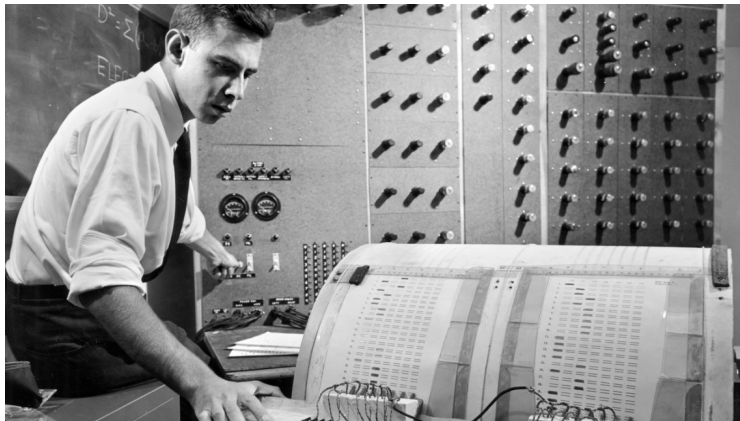
2015 VGG19 (Oxford): Convolutional, 1.45×10^8

2020 GPT-3 (OpenAI): Transformer, 1.75×10^{11}

Bernstein, Liane, et al. "Freely scalable and reconfigurable optical hardware for deep learning." Scientific reports 11.1 (2021): 3144.

Perceptron

F. Rosenblatt. The Perceptron: A Probabilistic Model For Information. Storage And Organization In The Brain. Psychological Review 65 (6): 386-408.



Fuente: <https://news.cornell.edu/stories/2019/09/professors-perceptron-paved-way-ai-60-years-too-soon>

Perceptron: modelo de una neurona

- Modelo **paramétrico**:

$$\hat{y} = f_{\theta}(\mathbf{x})$$

donde la función f está parametrizada por θ

- **Regresión lineal** - Modelo paramétrico **básico** de regresión

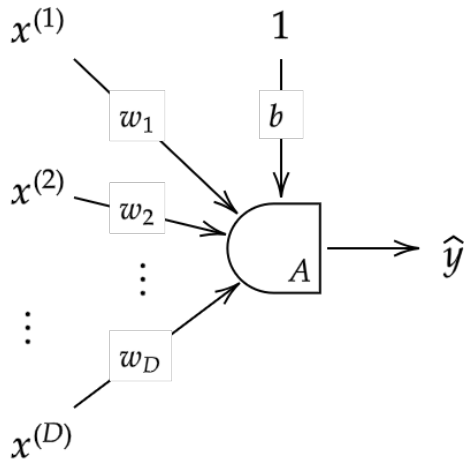
$$\hat{y} = w_1x^{(1)} + w_2x^{(2)} + \dots + w_Dx^{(D)} + b = \mathbf{x}^{\top} \mathbf{w} + b$$

- **Regresión logística** - Modelo paramétrico **básico** de clasificación

$$\hat{y} = \text{Sigmoid} (w_1x^{(1)} + w_2x^{(2)} + \dots + w_Dx^{(D)} + b) = \text{Sigmoid} (\mathbf{x}^{\top} \mathbf{w} + b)$$

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

Perceptron: modelo de una neurona



Usando la notación matricial:

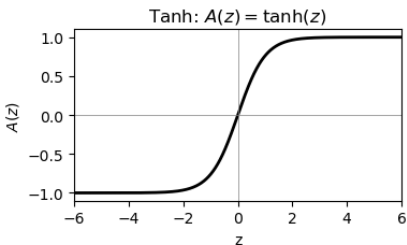
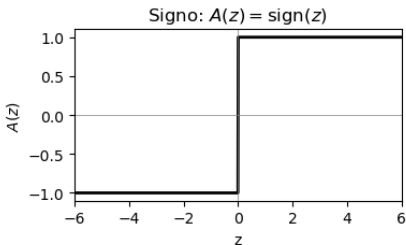
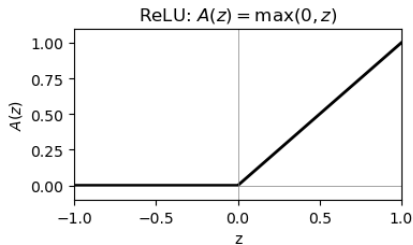
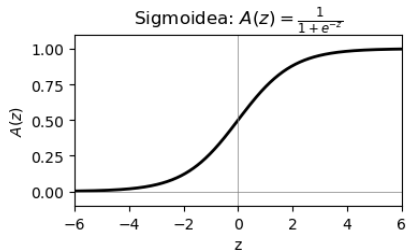
$$\hat{y} = A (\mathbf{x}^\top \mathbf{w} + b)$$

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



Componente básica: Densa

El perceptron es un caso particular de una **componente básica**, llamada **densa**:

$$\hat{\mathbf{y}} = \mathcal{D}(\mathbf{X}; \mathbf{W}, \mathbf{b}) = A(\mathbf{XW} + \mathbf{b}) \in \mathbb{R}^{N \times K}$$

con

$$\mathbf{X} \in \mathbb{R}^{N \times D}, \mathbf{y} \in \mathbb{R}^{N \times K}, \mathbf{W} \in \mathbb{R}^{D \times K}, \mathbf{b} \in \mathbb{R}^{N \times K}$$

- el producto $\mathbf{XW} + \mathbf{b}$ es la **unidad lineal**
- A es la función de **activación**

Multi-layer Perceptron

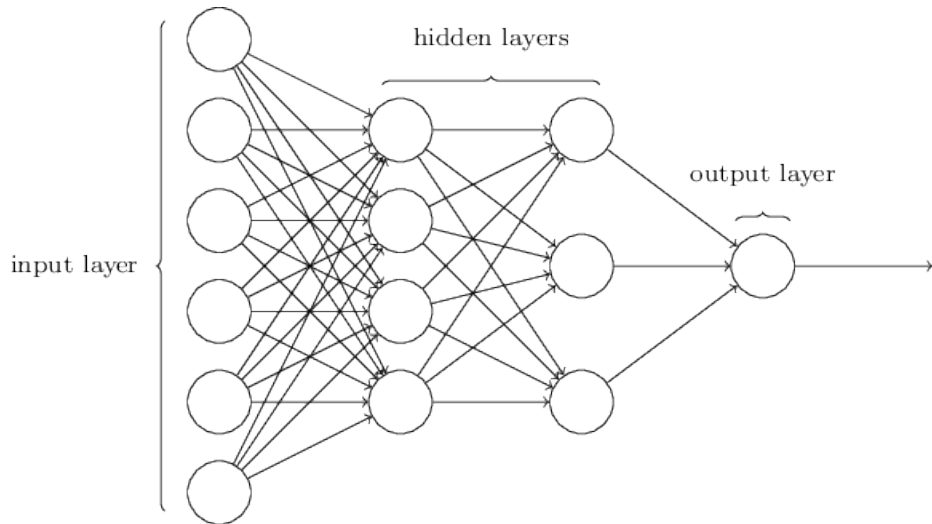
- Un **MLP** de $d > 1$ capas se obtiene componiendo d capas densas:

$$M(\mathbf{X}; \mathbf{W}_1, \dots, \mathbf{W}_d, \mathbf{b}_1, \dots, \mathbf{b}_d) = D_d(D_{d-1}(\dots, \mathbf{W}_{d-1}, \mathbf{b}_{d-1}); \mathbf{W}_d, \mathbf{b}_d)$$

- La dimensión de salida de D_i tiene que ser igual a la de entrada de D_{i+1}
- La cantidad de *parámetros* de M es:

$$\text{params}(M) = \sum_{i=1}^d \text{params}(D_i) = \sum_{i=1}^d D_i K_i + K_i \quad \mathbf{W}_i \in \mathbb{R}^{D_i \times K_i}, \mathbf{b}_i \in \mathbb{R}^{1 \times K_i}$$

MLP: ejemplo gráfico



Funciones de Pérdida

- **Regresión (MSE):**

$$L(\hat{y}, y) = (\hat{y} - y)^2$$

- **Clasificación binaria (BCE):** la última función de activación es una **sigmoidea**

$$L(\hat{y}, y) = -\left(y \log \hat{y} + (1 - y) \log(1 - \hat{y})\right)$$

- **Clasificación con K clases (CCE):** la última función de activación es una **softmax** e $y \in \{0, 1\}^K$ (one-hot)

$$L(\hat{y}, y) = -\sum_k y_k \log \hat{y}_k$$

Costo Empírico

- Debemos elegir una función de **pérdida** $L(\hat{y}, y)$
- Para un (\mathbf{W}, \mathbf{b}) dado, la **costo empírico** es:

$$J(\mathbf{W}, \mathbf{b}) = \frac{1}{N} \sum_{i=1}^N L(M(\mathbf{x}_i; \mathbf{W}, \mathbf{b}), y_i)$$

donde $M(\mathbf{x}_i; \mathbf{W}, \mathbf{b})$ es la salida de la red para la entrada \mathbf{x}_i .

- **Entrenar** la red significa encontrar el argumento mínimo de J (**ERM**):

$$(\hat{\mathbf{W}}, \hat{\mathbf{b}}) = \arg \min_{(\mathbf{W}, \mathbf{b})} J(\mathbf{W}, \mathbf{b}).$$

ERM: entrenamiento del MLP

- Problema de optimización - variantes de **descenso por gradiente**:

$(\mathbf{W}_0, \mathbf{b}_0)$ = inicializar random

$$(\mathbf{W}_{k+1}, \mathbf{b}_{k+1}) = (\mathbf{W}_k, \mathbf{b}_k) - \eta \cdot \nabla J(\mathbf{W}_k, \mathbf{b}_k)$$

donde η es la **tasa de aprendizaje**

- Debemos **calcular** el **gradiente de J** en cada paso.

- **Linealidad**:

$$\nabla J(\mathbf{W}, \mathbf{b}) = \frac{1}{N} \sum_{i=1}^N \nabla L(M(\mathbf{x}_i; \mathbf{W}, \mathbf{b}), y_i) = \mathbf{E}_{(\mathbf{x}, y) \sim T} [\nabla L(M(\mathbf{x}; \mathbf{W}, \mathbf{b}), y)]$$

ERM: entrenamiento del MLP

Algoritmo básico: Descenso del gradiente estocástico (SGD)

- Comenzar con $(\mathbf{W}_0, \mathbf{b}_0)$ random

- En cada **step** t : seleccionar $(\mathbf{x}, y) \in \mathcal{T}$ aleatoriamente y actualizar

$$(\mathbf{W}_{t+1}, \mathbf{b}_{t+1}) = (\mathbf{W}_t, \mathbf{b}_t) - \eta \nabla L(M(\mathbf{x}; \mathbf{W}_t, \mathbf{b}_t), y) \quad \eta \text{ es el learning rate}$$

En la práctica se actualiza por **batch** (se aprovecha el paralelismo)

Se usa una técnica llamada **back-propagation** basada en la regla de derivación de la cadena $d(f \circ g) = df \cdot dg$ para entrenar un MLP de muchas capas

1986

Rumelhart, Hinton & Williams introduce **Backpropagation Algorithm**

Learning representations by back-propagating errors

**David E. Rumelhart*, Geoffrey E. Hinton†
& Ronald J. Williams***

* Institute for Cognitive Science, C-015, University of California,
San Diego, La Jolla, California 92093, USA

† Department of Computer Science, Carnegie-Mellon University,
Pittsburgh, Philadelphia 15213, USA



We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight adjustments, internal 'hidden' units which are not part of the input or output come to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-convergence procedure¹.

Nature. 323 (6088): 533-536.

Backpropagation (Opcional)

Objetivo: Calcular $\nabla L(\hat{y}, y)$ eficientemente para un MLP con K capas.

- $\mathbf{z}_k = \mathbf{W}_k \mathbf{a}_{k-1} + \mathbf{b}_k$ (pre-activación)

- $\mathbf{a}_k = \sigma(\mathbf{z}_k)$ (post-activación)

- $\mathbf{a}^0 = \mathbf{x}$

Forward pass:

- Para $k = 1$ hasta K : calcular \mathbf{z}_k y \mathbf{a}_k

Backward pass (errores):

- Calcular $\delta_K = \sigma'(\mathbf{z}_K) \nabla_{\mathbf{a}_K} L$

- Para $k = K - 1$ hasta 1: $\delta_k = \sigma'(\mathbf{z}_k) \mathbf{W}_{k+1}^\top \delta_{k+1}$

Backpropagation (Opcional)

Gradientes:

Para cada capa $k = 1, \dots, K$:

$$\nabla_{\mathbf{w}_k} L = \delta_k \cdot \mathbf{a}_{k-1}^\top$$

$$\nabla_{\mathbf{b}_k} L = \delta_k$$

Actualización (ej. SGD):

$$\mathbf{w}_k \leftarrow \mathbf{w}_k - \eta \nabla_{\mathbf{w}_k} L, \quad \mathbf{b}_k \leftarrow \mathbf{b}_k - \eta \nabla_{\mathbf{b}_k} L$$

Regularización

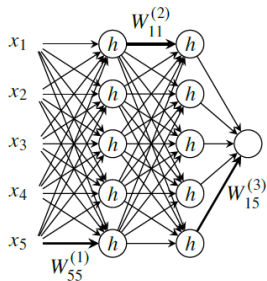
Se usan diferentes **regularizadores**:

- **Ridge** (norma l_2), **lasso** (norma l_1)
- **Dropout**: no actualizar ciertos coeficientes seleccionados aleatoriamente (efecto “ensemble”)
- **Early stopping**: usar un conjunto de validación V y parar de entrenar cuando el error en V no decrece por un cierto tiempo
- **Data augmentation**: generar datos a través de transformaciones (típico en análisis de imágenes: rotación, zoom, ...)

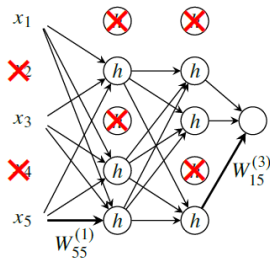
Dropout

Entrenamiento

Máscara que apaga aleatoriamente neuronas de una capa con probabilidad $1 - r$



(a) A standard neural network

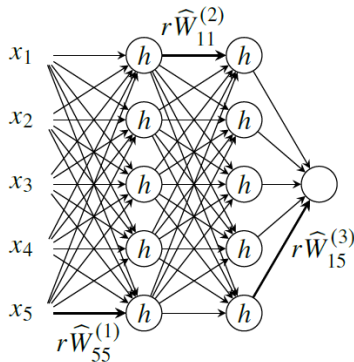


(b) Two sub-networks

Dropout

Predicción

Se multiplican los pesos por la probabilidad r



Early Stopping

Early Stopping

- 01: Inicializar $W_{\text{best}} \leftarrow \text{random}$
- 02: Inicializar $L_{\text{best}} \leftarrow +\infty$
- 03: Inicializar p paciencia
- 04: Inicializar $c \leftarrow 0$ contador paciencia
- 05:
- 06: Para cada epoch en epochs:
- 07: Entrenar red con datos de entrenamiento
- 08: $L \leftarrow$ evaluar red con datos de validación
- 09: Si $L < L_{\text{best}}$:
- 10: $L_{\text{best}} \leftarrow L$
- 11: $W_{\text{best}} \leftarrow$ pesos actuales de la red
- 12: $c \leftarrow 0$
- 13: Sino:
- 14: $c \leftarrow c + 1$
- 15: si $c \geq p$: break
- 16: Establecer pesos de la red $\leftarrow W_{\text{best}}$

Bibliografía

- Deep Learning. Ian Goodfellow, Yoshua Bengio, Aaron Courville. MIT Press, 2016.
- Prince, Simon JD. Understanding deep learning. MIT press, 2023.
- Bishop, Christopher M., and Hugh Bishop. Deep learning: Foundations and concepts. Springer Nature, 2023.
- Scardapane, Simone. "Alice's Adventures in a Differentiable Wonderland—Volume I, A Tour of the Land." arXiv preprint arXiv:2404.17625 (2024).