

# Cursos Extraordinarios

## Verano 2024

**“Inteligencia Artificial y Grandes Modelos de Lenguaje: Funcionamiento, Componentes Clave y Aplicaciones”**

**Zaragoza, del 3 al 5 de julio**

# Introducción al Aprendizaje Automático y las Redes Neuronales

# Deep Learning

---

- Rama de Machine Learning:
  - Aproximación de propósito general que permite aprender a partir de datos.
  - Basada en el uso de Redes Neuronales Artificiales
- Asociado actualmente al concepto de Inteligencia Artificial

**Capacidad de las máquinas para realizar tareas que normalmente requieren destrezas humanas.**

**Esto incluye:**

**Aprendizaje, percepción, razonamiento, reconocimiento de voz, comprensión del lenguaje natural, toma de decisiones, ...**



# Redes Neuronales Artificiales

- **Antecedentes:**

*La teoría:*

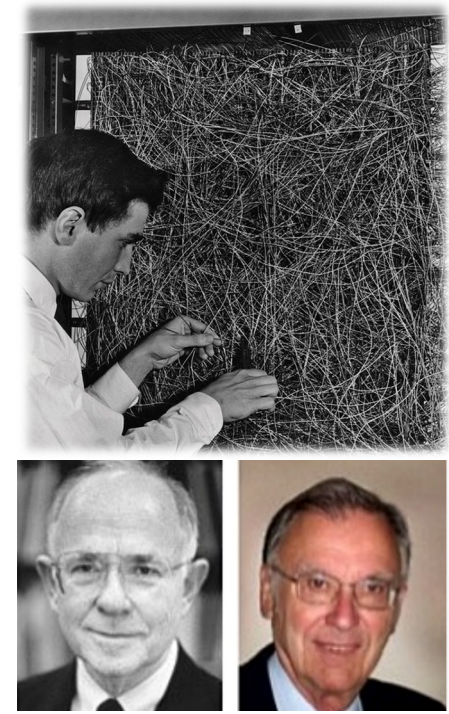
- McCulloch & Pitts 1943
- Wiener 1948
- Alan Turing 1950
- Frank Rosenblatt 1957



*Máquinas electrónicas:*

- Mark I, 1944
- Widrow & Hoff 1960

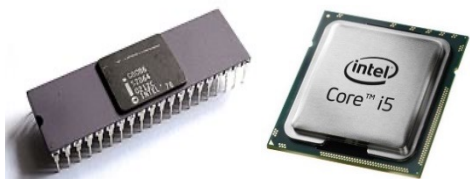
1 operación cada 3 segundos



# Explosión Tecnológica

## Microprocesadores

Intel 8086, 1978  
 50 mil operaciones por segundo  
 Intel i5, 2018  
 25 mil millones de operaciones por segundo



### 2010s La era de las GPUs

Playstation 4s, 2016  
 1.8 TFlops (~90 x intel i5)  
 Playstation 5s, 2020  
 10.2 TFlops (~411 x intel i5)

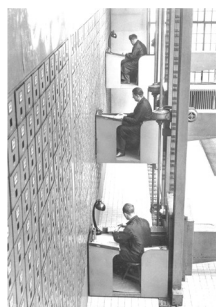


Nvidia RTX Titan, 2018  
 16 TFlops (~640 x intel i5)  
 Nvidia RTX 3090, 2020  
 35 TFlops (~1400 x intel i5)  
 Nvidia RTX 4090, 2022  
 82 TFlops (~3280 x intel i5)



## Almacenamiento

Sistema mecánico 1937 (República Checa)



### Capacidad de almacenamiento

|                 |      |        |
|-----------------|------|--------|
| Cinta perforada | 1970 | <1 KB  |
| Disco 3 1/2     | 1987 | 1.4 MB |
| DVD             | 1995 | 4.7 GB |



### Velocidad de almacenamiento

|                  |      |                 |
|------------------|------|-----------------|
| Disco duro       | 2000 | 18GB (48MB/s)   |
| HD estado sólido | 2021 | 1TB (7000 MB/s) |



## Comunicaciones / Difusión

Proceedings papel 1995 / Revistas papel



Buscadores internet 1998



Software gratuito / Toolkits 2010

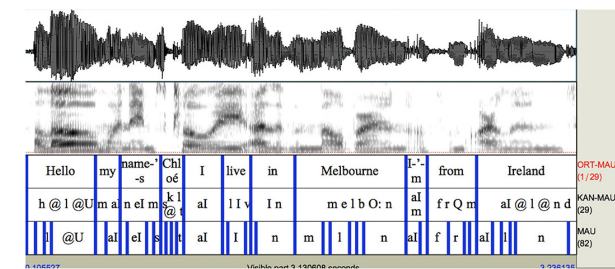
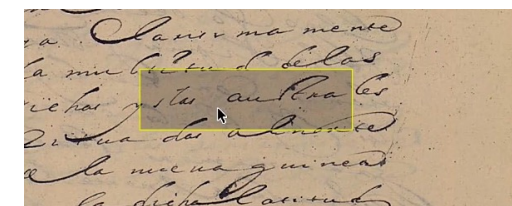
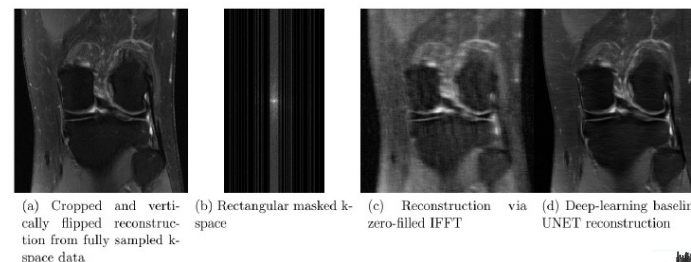
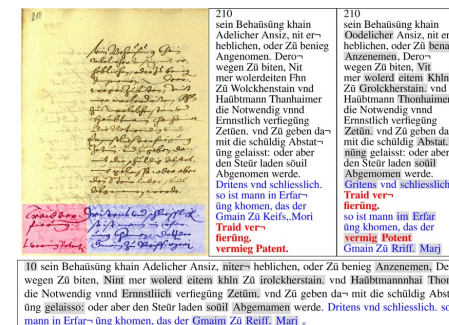


2008 Redes sociales / plataformas de desarrollo colaborativo



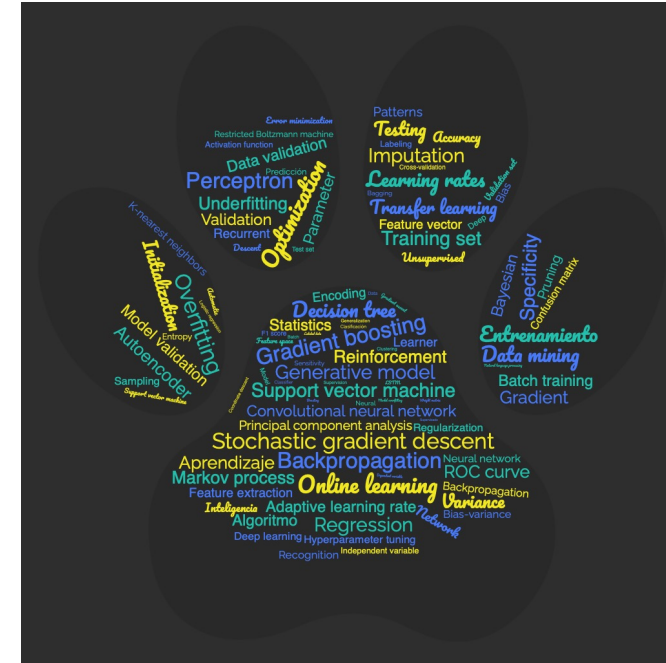
# Campos de Aplicación

- Visión por Computador
- Reconocimiento de Escritura
- Ayuda al diagnóstico
- Transcripción de voz
- Traducción
- Algoritmos de Recomendación
- ...



# Aprendizaje Automático

- “*Learn from data*”
- Encontrar patrones y tendencias, entendiendo “lo que nos dicen los datos”
- Sistema automático que aprende de la **Experiencia (E)** para realizar una **Tarea (T)** determinada evaluada a través de una **Métrica (M)** establecida.

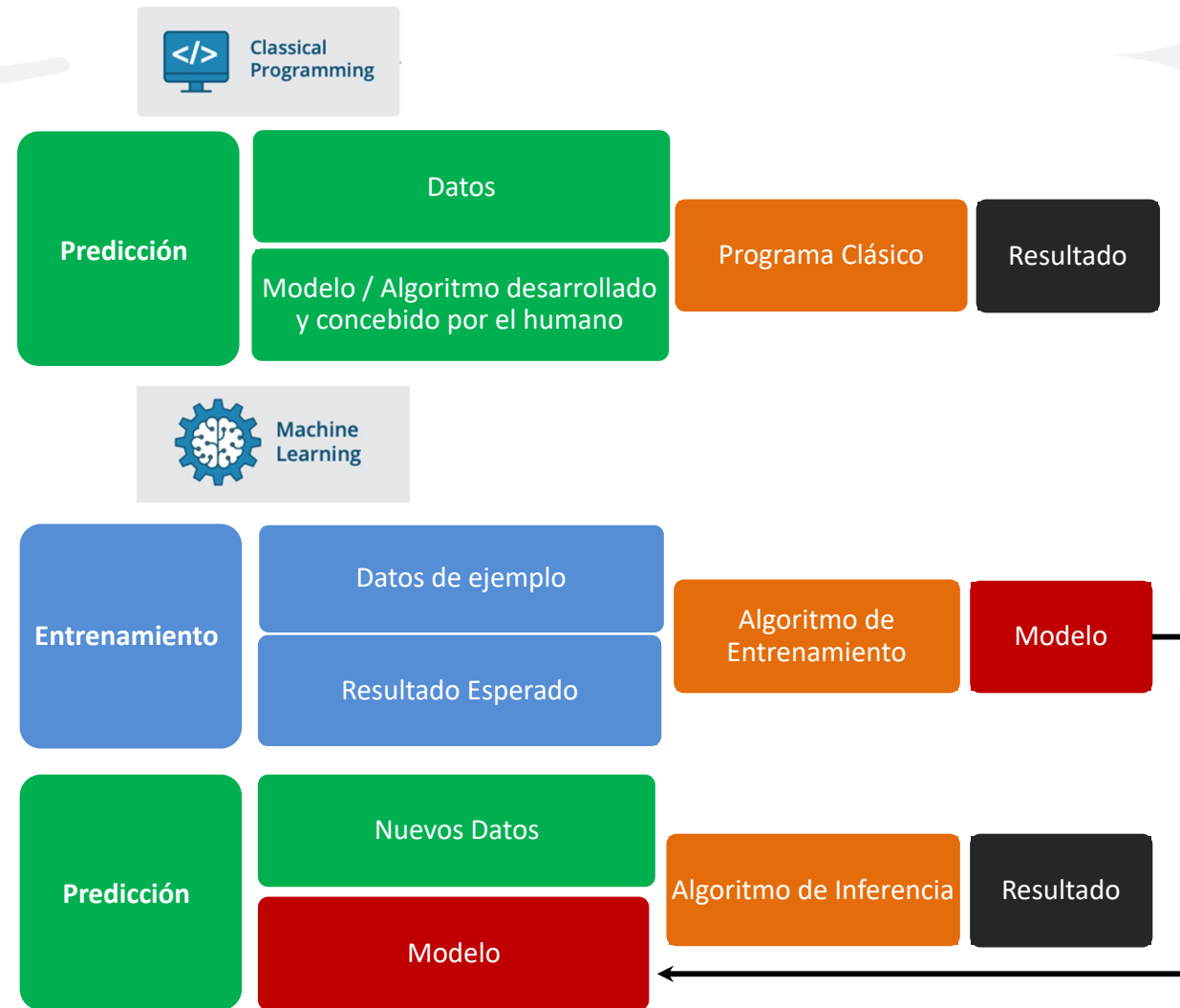


**Tom Michael Mitchell**

American computer scientist; Founders University Professor  
[Carnegie Mellon University](https://www.cmu.edu/about/carnegie-mellon-university) (CMU).

Tom Mitchell [[T. Mitchell. Machine Learning. McGraw Hill, 1997](#)]

# Aprendizaje Automático

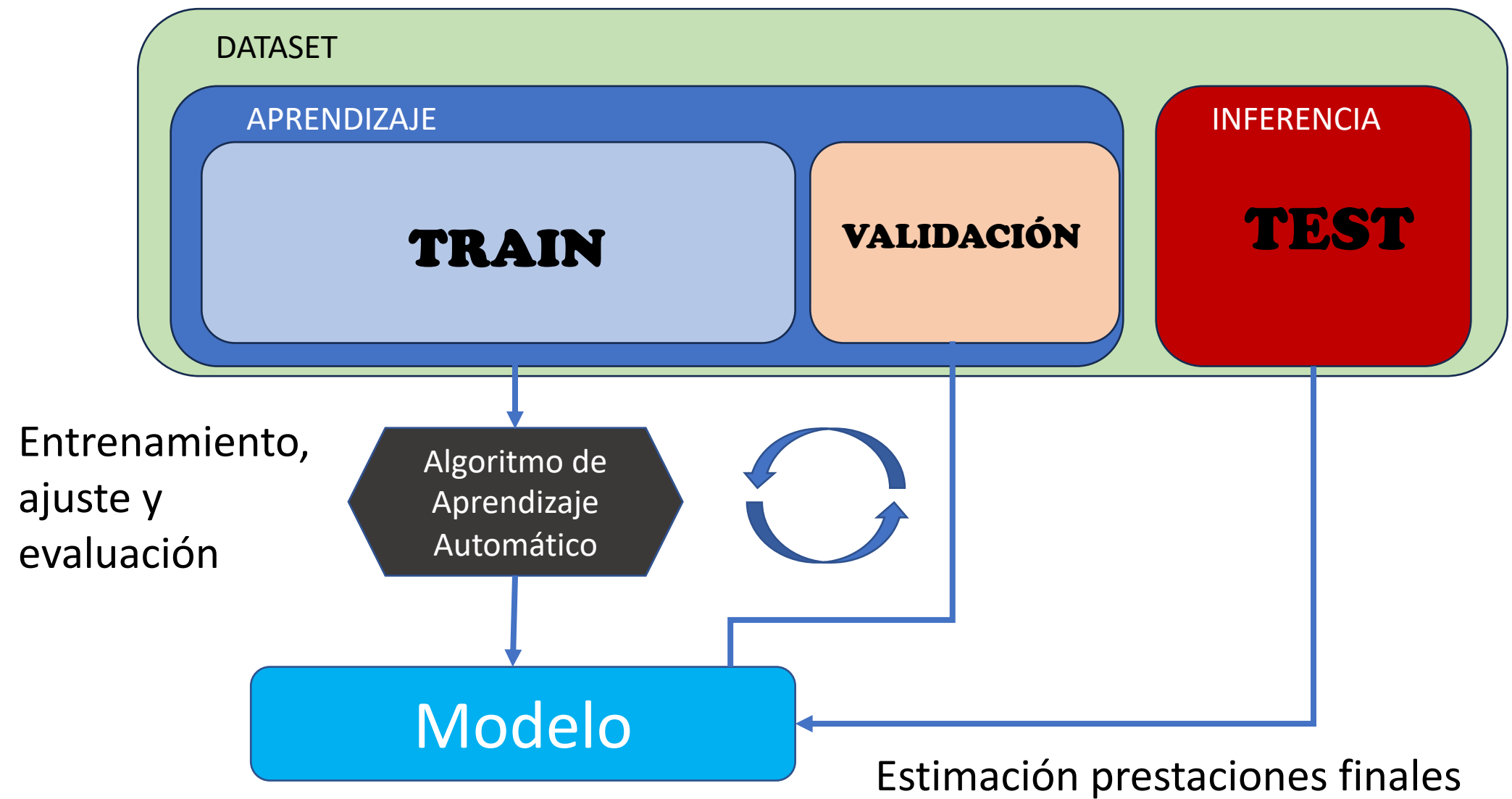


Ofrece soluciones a **Problemas** demasiado **complejos** como para ser resueltos por **programas clásicos** concebidos por humanos



# Diseño Experimental

- **Uso de los datos:**



# Tipos de Tareas

---

- **Clasificación:**
  - Simple:
  - Múltiple:
- **Regresión:**
  - Simple:
  - Múltiple:

# Tipos de Tareas

- **Clasificación:**

## CLASIFICACIÓN SIMPLE

- Simple:

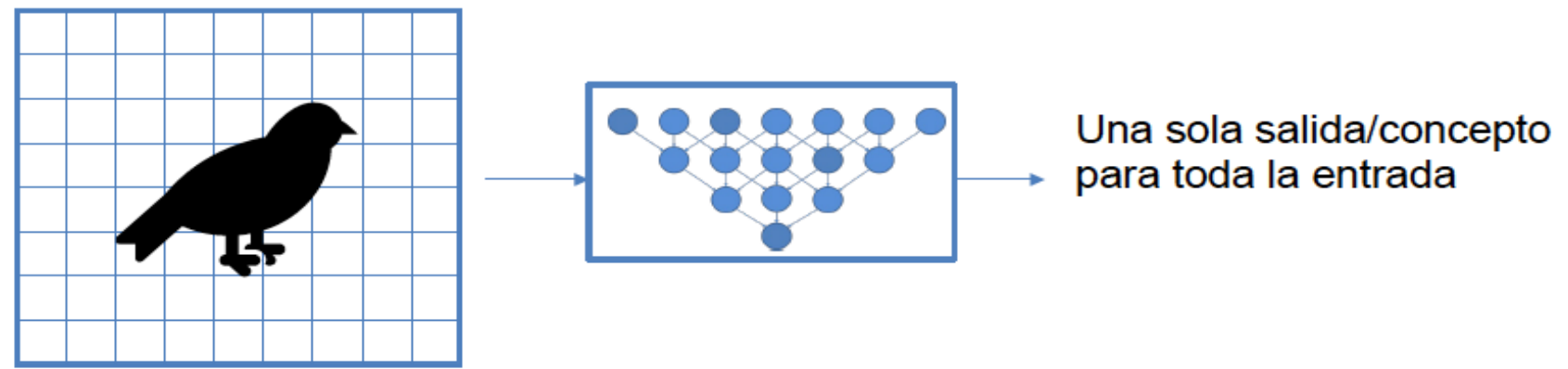
**¿Qué CONCEPTO hay tras una imagen/texto/audio ?**

- Múltiple:

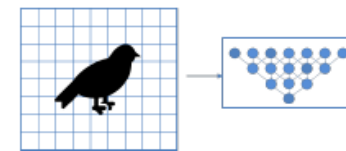
- **Regresión:**

- Simple:

- Múltiple:



# Tipos de Tareas



- Clasificación:**

## CLASIFICACIÓN SIMPLE

- Simple:

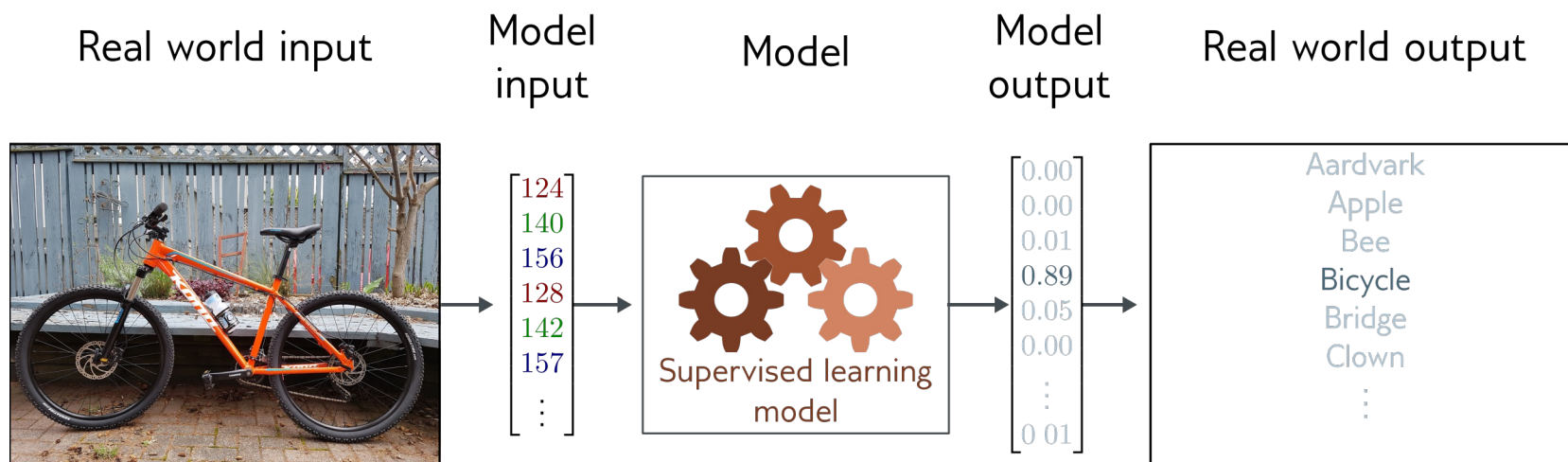
**¿Qué CONCEPTO hay tras una imagen/texto/audio ?**

- Múltiple:

- Regresión:**

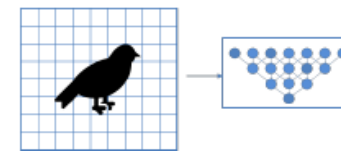
- Simple:

- Múltiple:



**A cada entrada, el modelo decide UNA entre C posibles respuestas (Clases)**

# Tipos de Tareas



- **Clasificación:**

## CLASIFICACIÓN SIMPLE

- Simple:

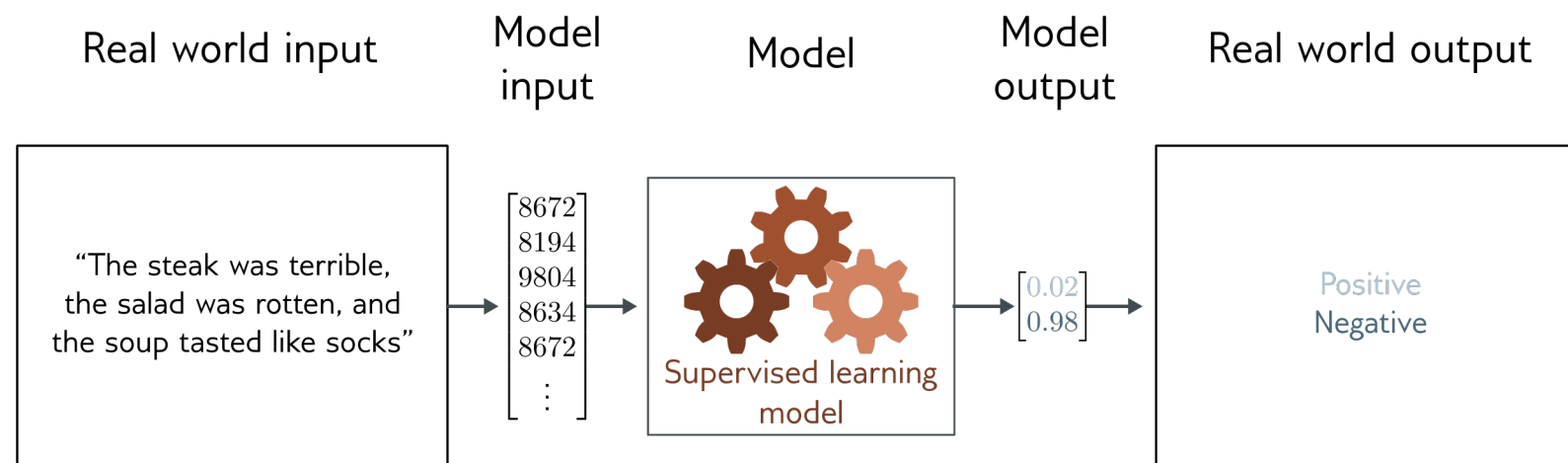
- Múltiple:

¿Qué **CONCEPTO** hay tras una imagen/texto/audio ?

- **Regresión:**

- Simple:

- Múltiple:



En muchas ocasiones **C=2** (Caso binario SI/NO)

# Tipos de Tareas

- Clasificación:**

## CLASIFICACIÓN MÚLTIPLE

- Simple:

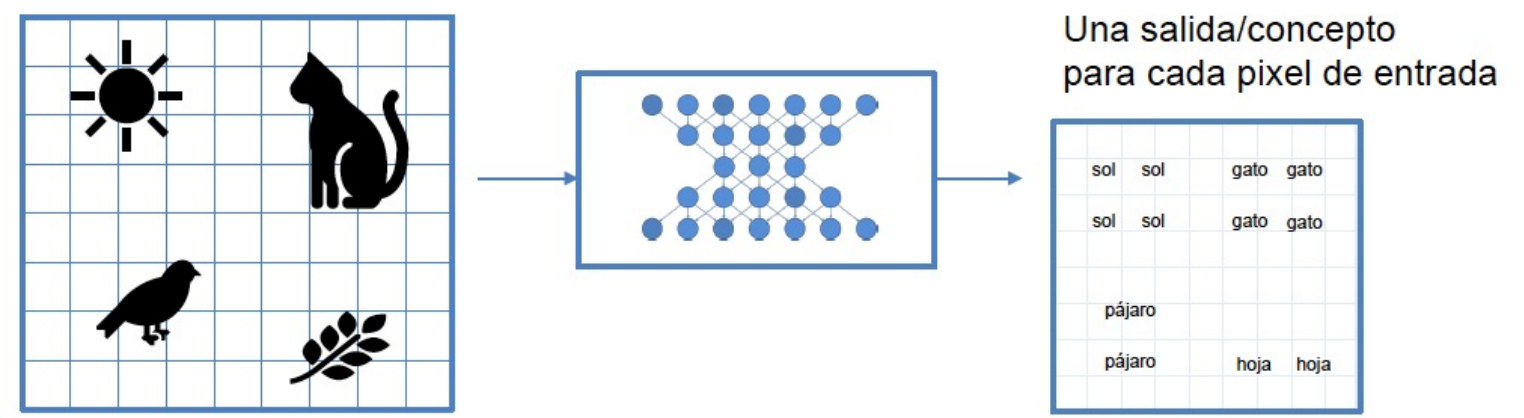
- Múltiple:

**¿Qué CONCEPTO hay en cada zona/pixel/fragmento de una imagen/texto/audio ?**

- Regresión:**

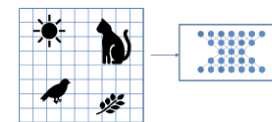
- Simple:

- Múltiple:



**Enunciar VARIAS PROPIEDADES/CONCEPTOS de una zona/pixel/fragmento de una imagen/texto/audio ?**

# Tipos de Tareas



- **Clasificación:**

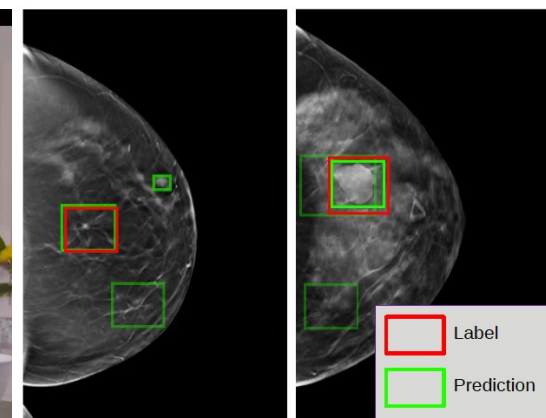
## CLASIFICACIÓN MÚLTIPLE

- Simple: **¿Qué CONCEPTO hay en cada zona/píxel/fragmento de una imagen/texto/audio ?**

- **Múltiple:**

- **Regresión:**

- Simple:
- Múltiple:



**DBTex Challenge:**  
SPIE-AAPM-NCI DAIR Digital Breast  
Tomosynthesis Lesion Detection Challenge



**Enunciar VARIAS PROPIEDADES/CONCEPTOS de una zona/píxel/fragmento de una imagen/texto/audio ?**

# Tipos de Tareas

- **Clasificación:**

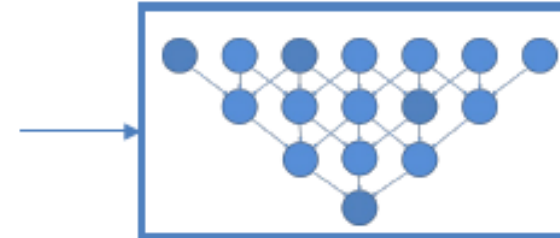
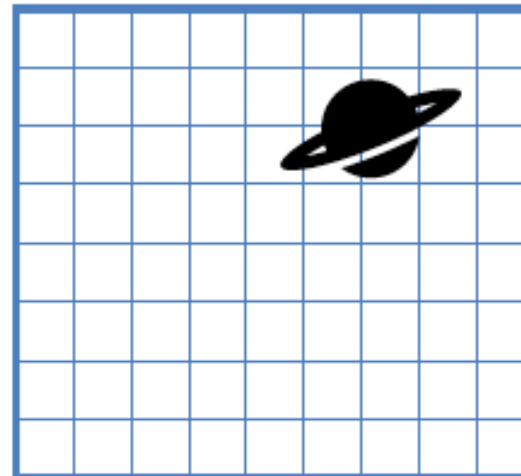
- Simple:
- Múltiple:

## REGRESIÓN SIMPLE

**Realizar una PREDICCIÓN NUMÉRICA a partir de una imagen/texto/audio ?**

- **Regresión:**

- Simple:
- Múltiple:



Estimar la posición  
Coordenadas  
 $X = 8.7$   
 $Y = 6.6$

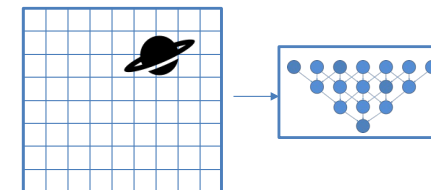


# Tipos de Tareas

- **Clasificación:**

- Simple:
- Múltiple:

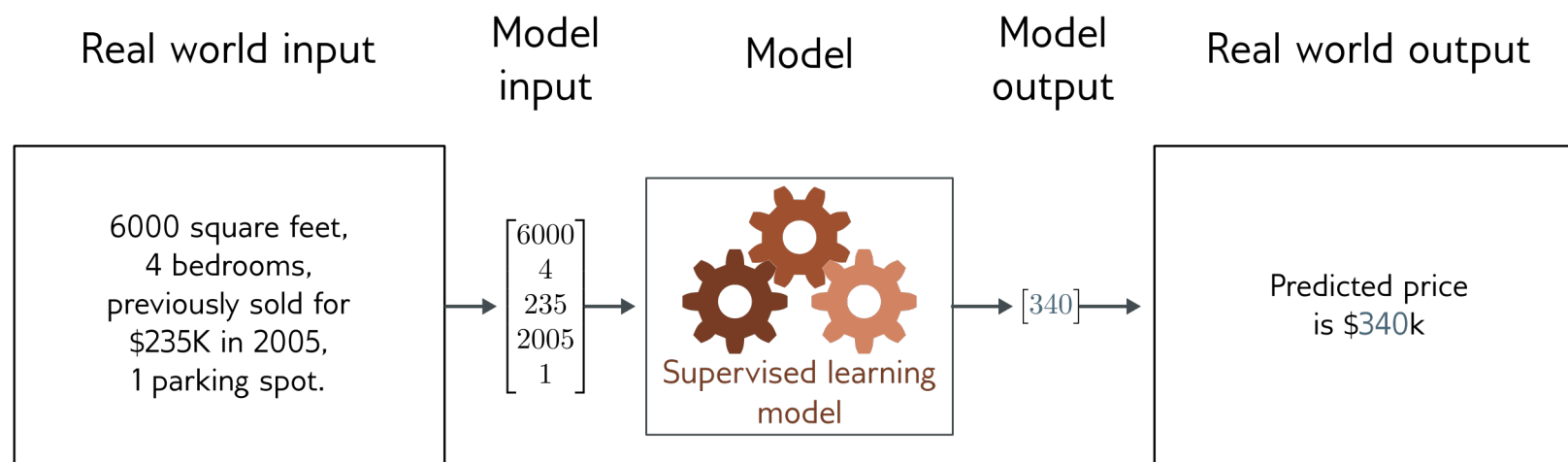
## REGRESIÓN SIMPLE



**Realizar una PREDICCIÓN NUMÉRICA a partir de una imagen/texto/audio ?**

- **Regresión:**

- Simple:
- Múltiple:



**Ante cada entrada, el modelo ofrece UNA respuesta**

# Tipos de Tareas

- **Clasificación:**

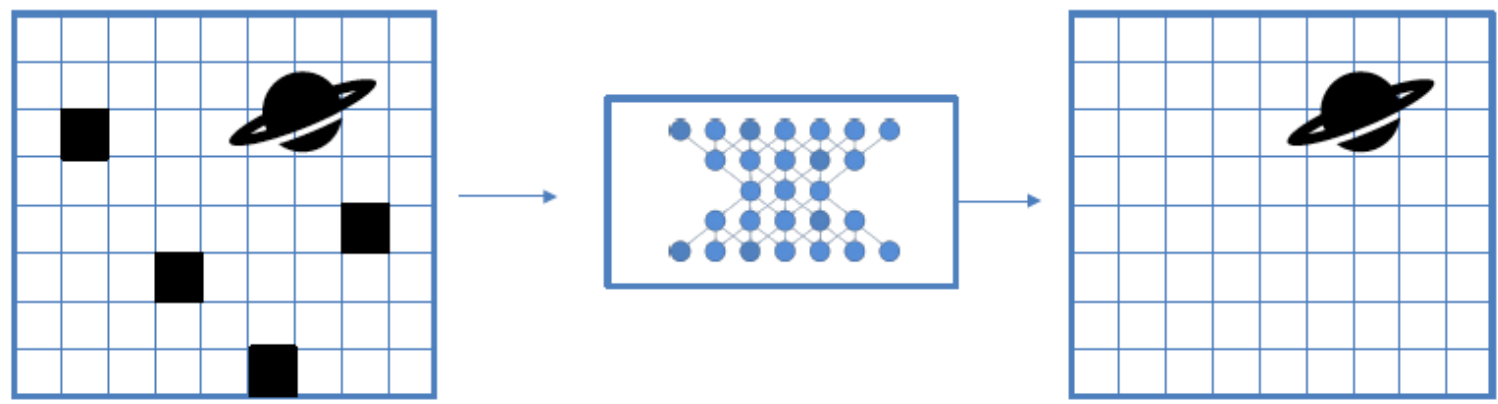
- Simple:
- Múltiple:

## REGRESIÓN MÚLTIPLE

**PREDECIR VARIOS VALORES NUMÉRICOS para cada zona/pixel/fragmento de una entrada**

- **Regresión:**

- Simple:
- Múltiple:



# MÉTRICAS DE EVALUACIÓN

- **Clasificación:**
  - Caso Binario

| <b>T = P+N</b>        |          | <b>VALOR REAL</b>  |  |
|-----------------------|----------|--|--|
|                       |          | POSITIVO (P = TP + FN)   | NEGATIVO (N = FN + TN)   |
| <b>VALOR PREDICHO</b> | POSITIVO | <b>TRUE POSITIVE (TP)</b>  | <b>FALSE POSITIVE (FN)</b><br><b>TYPE II ERROR</b><br><b>FALSA ALARMA</b><br><b>FALSA ACEPTACIÓN</b> |
|                       | NEGATIVO | <b>FALSE NEGATIVE (FN)</b><br><b>TYPE I ERROR</b><br><b>FALSO RECHAZO</b><br><b>MISS</b> | <b>TRUE NEGATIVE (TN)</b>  |

# MÉTRICAS DE EVALUACIÓN

## Clasificación:

- Caso Binario:

- Accuracy:**  $ACC = \frac{TP+TN}{P+N}$

|                |          | VALOR REAL   |  |
|----------------|----------|--|--|
|                |          | POSITIVO (P = TP + FN)                                       | NEGATIVO (N = FN + TN)   |
| VALOR PREDICHO | POSITIVO | TRUE POSITIVE (TP)   | FALSE POSITIVE (FN)<br>TYPE II ERROR<br>FALSA ALARMA<br>FALSA ACEPTACIÓN |
|                | NEGATIVO | FALSE NEGATIVE (FN)<br>TYPE I ERROR<br>FALSO RECHAZO<br>MISS | TRUE NEGATIVE (TN)   |

- Sensitivity / Recall / True Positive Rate:**  $TPR = \frac{TP}{P}$ 
  - También como complementario a **Miss rate/ False Negative Rate:**  $TPR = 1 - FNR$
- Especificidad / True Negative Rate:**  $TNR = \frac{TN}{N}$ 
  - También como complementario a **False Alarm Rate / False Positive Rate:**  $TNR = 1 - FPR$

POR COLUMNAS

- Valor Predictivo Positivo / Precisión:**  $PPV = \frac{TP}{TP+FP}$
- Valor Predictivo Negativo:**  $NPV = \frac{TN}{TN+FN}$

POR FILAS

# Tipos de Tareas

- **Clasificación:**

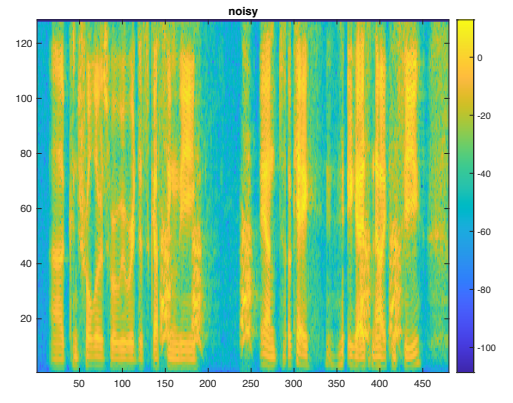
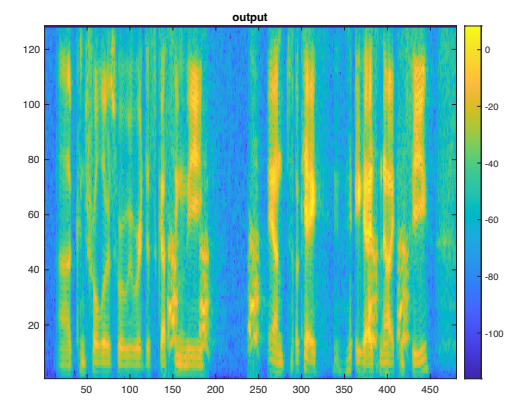
## REGRESIÓN MÚLTIPLE

- Simple:
- Múltiple:

**PREDECIR VARIOS VALORES NUMÉRICOS para cada zona/pixel/fragmento de una entrada**

- **Regresión:**

- Simple:
- Múltiple:

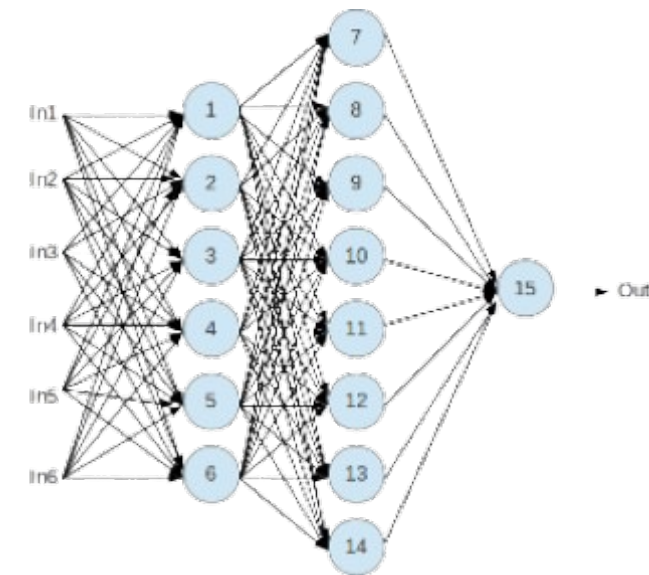


**Elimina el ruido de un audio ruidoso**



# Redes Neuronales Artificiales

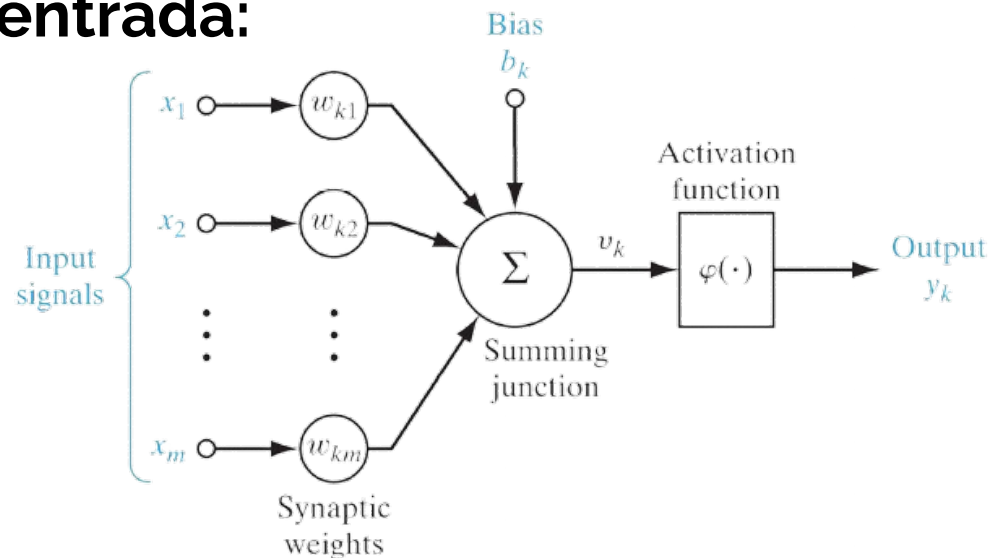
- **Conjunto de elementos simples de procesamiento interrelacionados:**
  - Distribución del procesamiento de la información
  - Cada unidad recibe información de otras, los agrega y transmite una respuesta a través de una función de activación
- **Cada elemento se conecta con otros**
  - Conexiones Sinápticas:
    - Cada conexión tiene un peso
    - Se ajustan a partir de ejemplos
      - Datos de entrenamiento
      - Algoritmos de aprendizaje



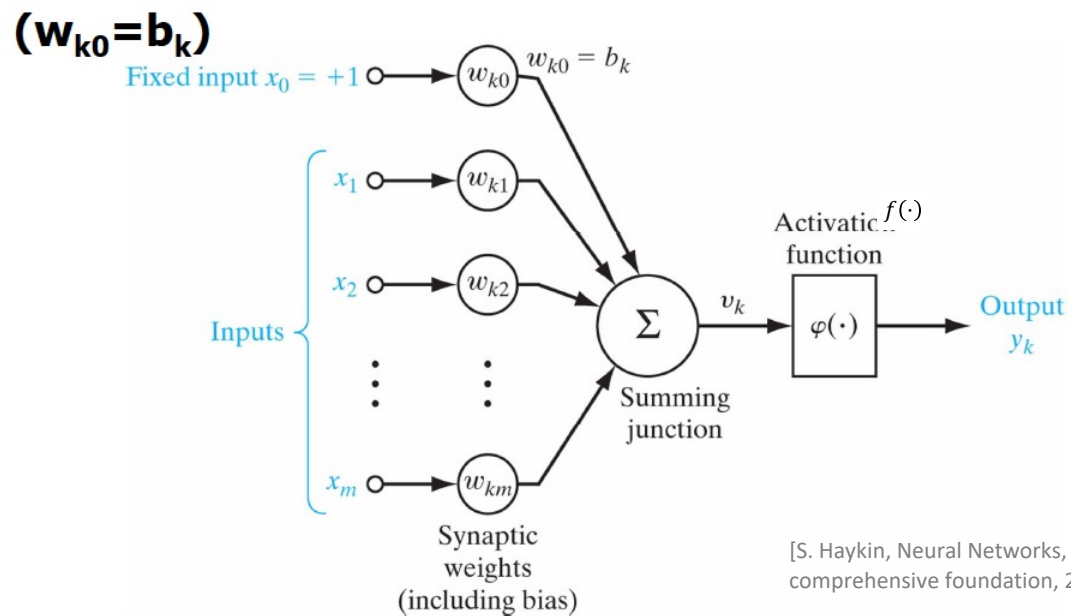
# Perceptrón

- Combinación lineal de valores a la entrada:

Vector de pesos  $\underline{w}$  + sesgo  $\underline{b}$



Transformación afín

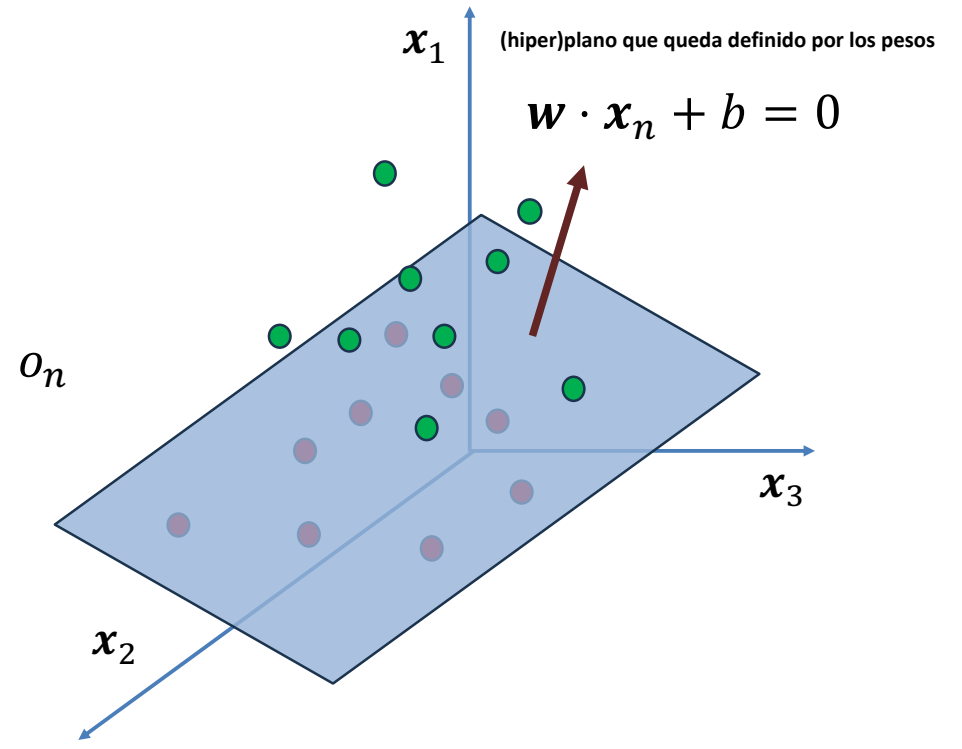
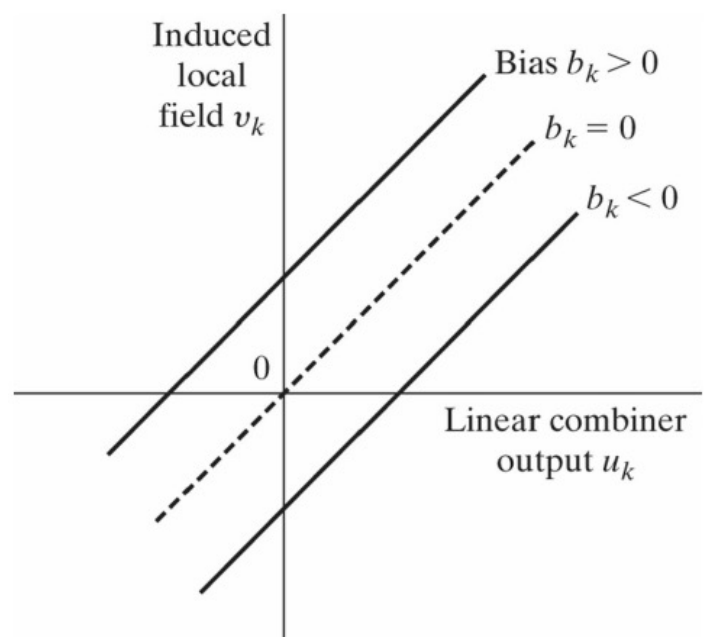


Función de activación,  $f(x_n)$ , no lineal

$$o_n = f(\mathbf{x}_n) = \begin{cases} +1, & \mathbf{w} \cdot \mathbf{x}_n + b > 0 \\ -1, & \mathbf{w} \cdot \mathbf{x}_n + b \leq 0 \end{cases}$$

[S. Haykin, Neural Networks, a comprehensive foundation, 2nd Edition]

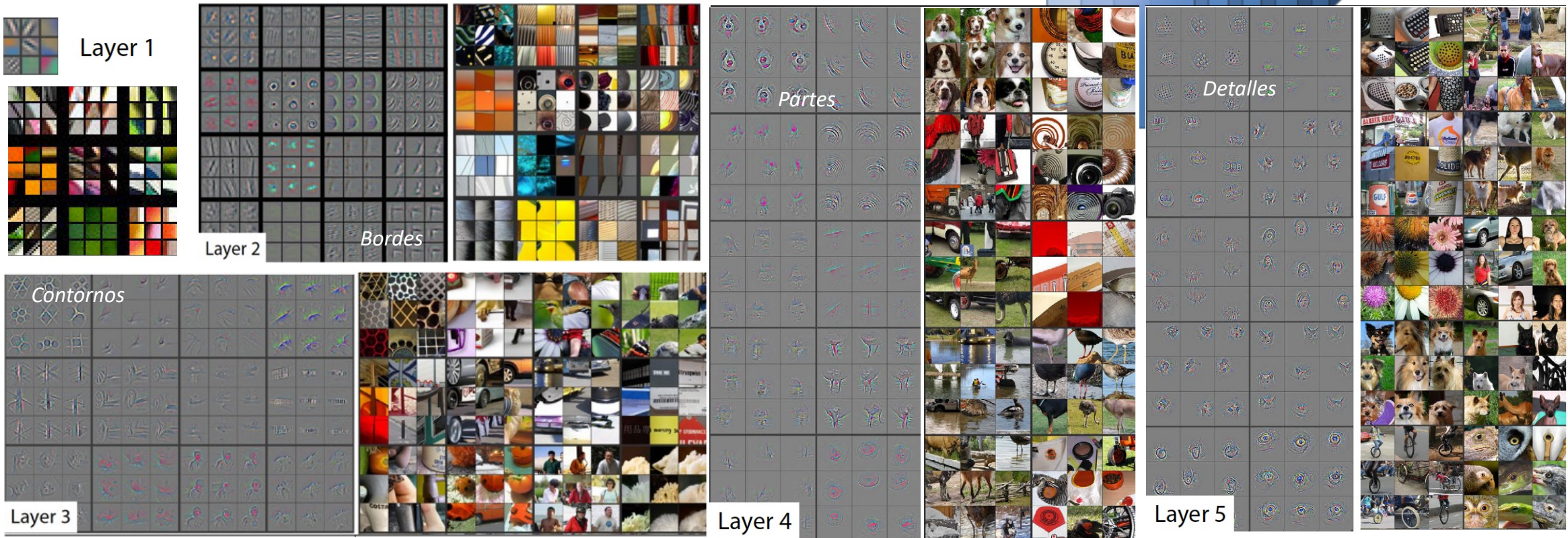
# Perceptrón





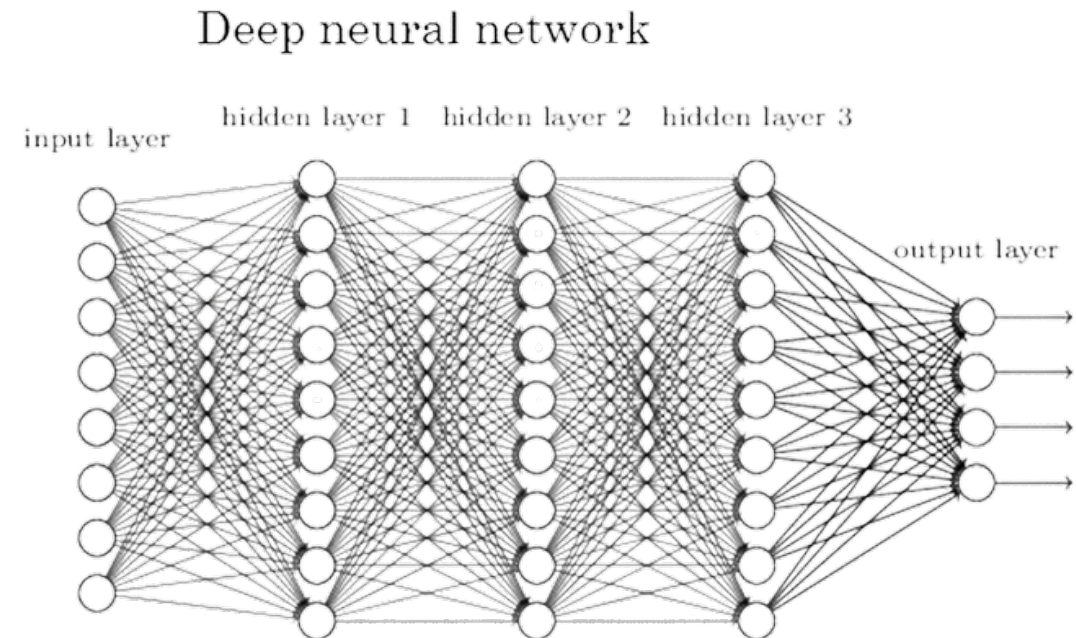
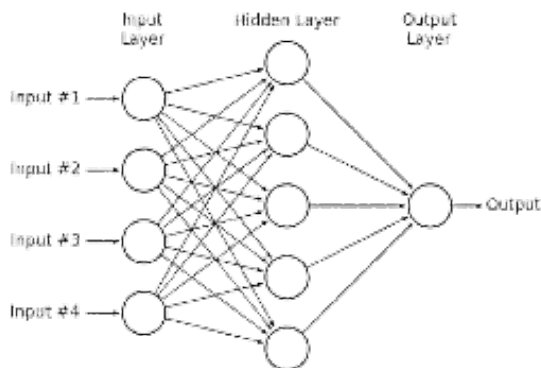
# Redes Neuronales Artificiales

- **Procesado en capas:**
  - Aprenden diferentes niveles de representación de los datos, desde lo más simple, hasta lo más abstracto



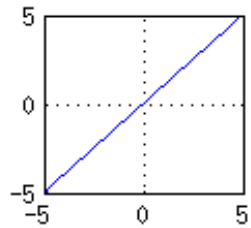
# Multilayer Perceptron / Deep Feedforward Network

- **Múltiples Neuronas (perceptrón) básicas organizadas en capas**
  - **Arquitectura de la red:**
    - **Capa de entrada:** Toma los datos de entrada (en forma vector)
      - Conviene acondicionar los datos: Normalización
    - **Capas Ocultas:** Sin contacto ni con la entrada ni con la salida
      - Agrupación de perceptrones
    - **Capas de Salida:**
      - Resultado del proceso



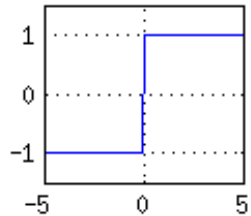


# Funciones de Activación



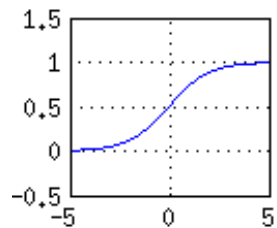
Linear

$$f(x) = x$$



Sign

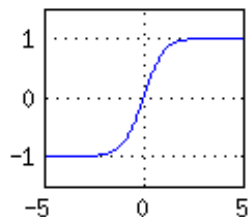
$$f(x) = \begin{cases} +1, & x \geq 0 \\ -1, & x < 0 \end{cases}$$



Logistic

$$f(x) = \frac{1}{1 + \exp(-x)}$$

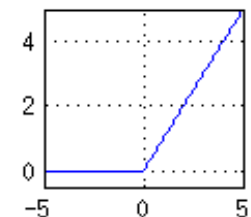
$$f'(x) = f(x)(1 - f(x))$$



Tanh

$$f(x) = \tanh(x)$$

$$f'(x) = (1 - f(x)^2)$$



ReLu

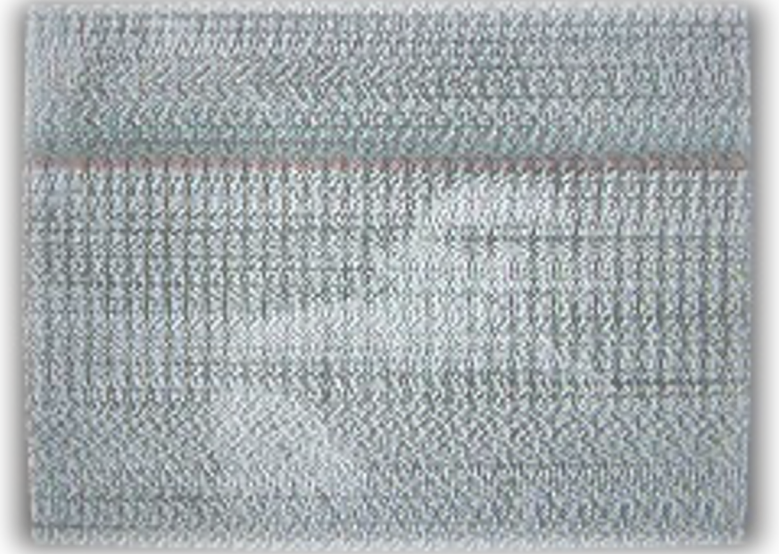
$$f(x) = \max(0, x)$$

$$f'(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

# Entrenamiento / Aprendizaje

---

- **Ajuste de los parámetros del modelo:**
  - Selección de los datos de entrenamiento
  - Exposición de los datos a la red
  - Observación de la salida
    - Variar los pesos para hacer disminuir el error



# Deep Forward Networks

---

- **Tipos de datos:**
  - **Etiquetas:**
    - Son vectores objetivo:  $\mathbf{y}_n$
    - Si la red actúa como regresor (prediciendo valores):
      - $\mathbf{y}_n \in \mathbb{R}^D$
    - Si la red actúa como un clasificador
      - $\mathbf{y}_n \in \{0, 1\}^C$  (C clases) One-hot encoding
      - $\mathbf{y}_n \in \{0, 1, \dots, C - 1\}$  (C clases)
  - **Salidas:**
    - Son el producto de la red:  $\mathbf{y}_n$
    - Tendrán la misma dimensión y serán del mismo tipo que las etiquetas
  - **Entradas:**
    - Son los datos que alimentan la red:  $\mathbf{x}_n$

# Deep Forward Networks

- Entrenamiento:**

- Pares de datos**  $(x_n, y_n)$ ,  $x_n \in \mathbb{R}^D$ ;  $y_n \in \mathbb{R}^D$  ó  $y_n \in \{0, 1\}^C \dots$

- Función de coste**  $J(X, Y, \Theta)$ :

- Es una función de los ejemplos  $X$ , de las etiquetas  $Y$  y de los parámetros del modelo  $\Theta$

- Error cuadrático:**

- $J = \sum_n (y_n - o_n)^2$

- Errores de Clasificación**

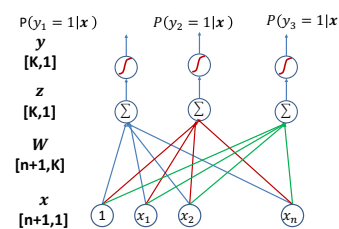
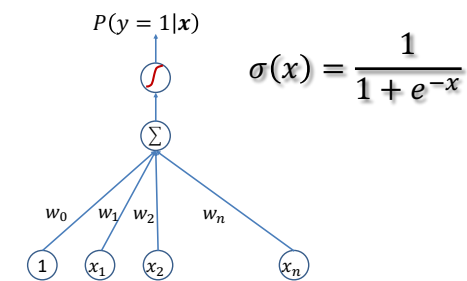
- $J(\hat{y}, y) = \sum_n L_{CE}(\hat{y}_n, y_n)$

Caso Binario

$$L_{CE}(\hat{y}_n, y_n) = -y_n \log y_n - (1 - y_n) \log(1 - \hat{y}_n) = -\log p(y_n | x_n)$$

Caso Multiclase

$$L_{CE}(\hat{y}_n, y_n) = -\sum_{k=1}^K y_n(k) \log \hat{y}_n(k)$$



$$\text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^K \exp(z_j)}$$

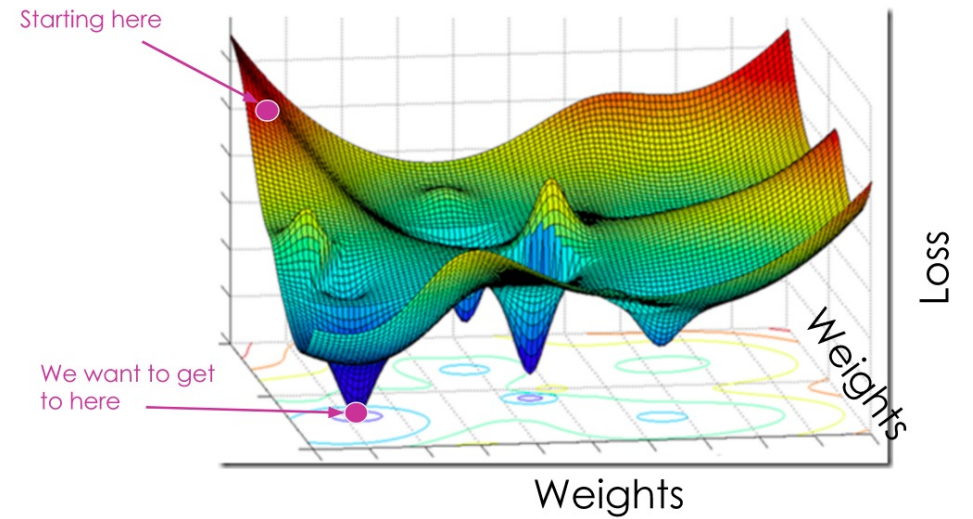
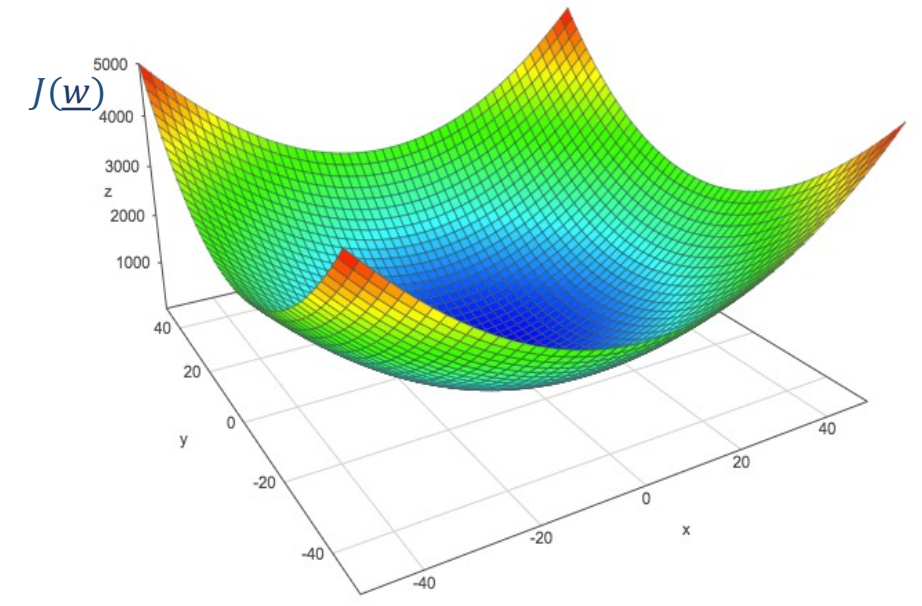
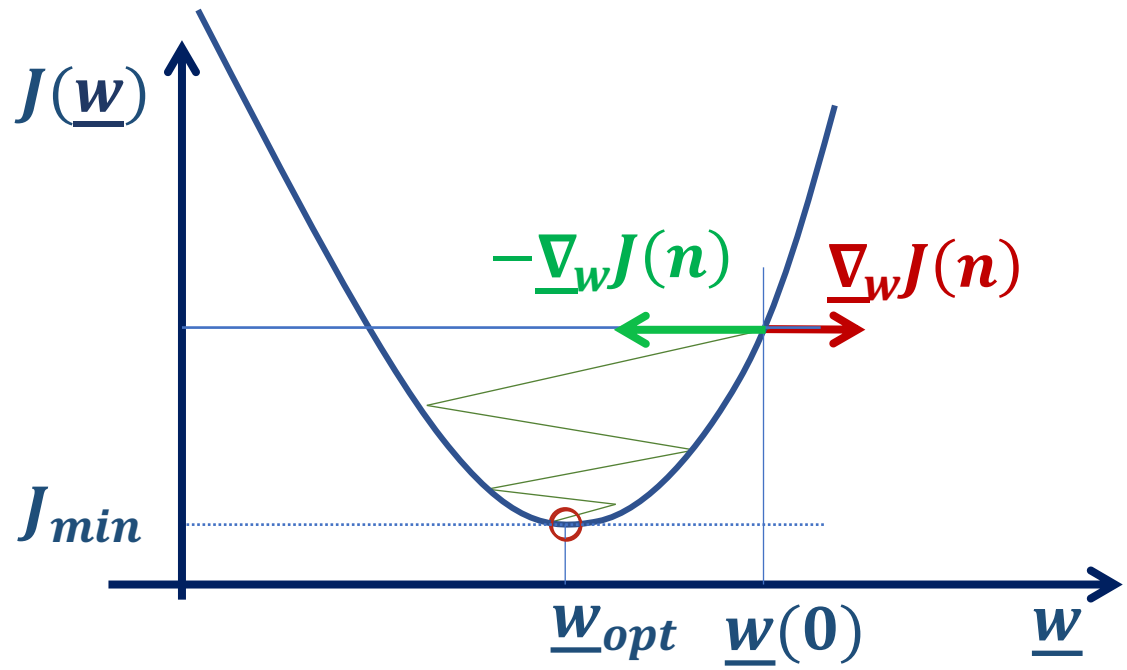
- Objetivo:**

- Encontrar los parámetros del modelo que minimicen la función de coste:

$$\{\hat{w}, \hat{b}\} = \arg \min_{w, b} J(X, Y, w, b)$$

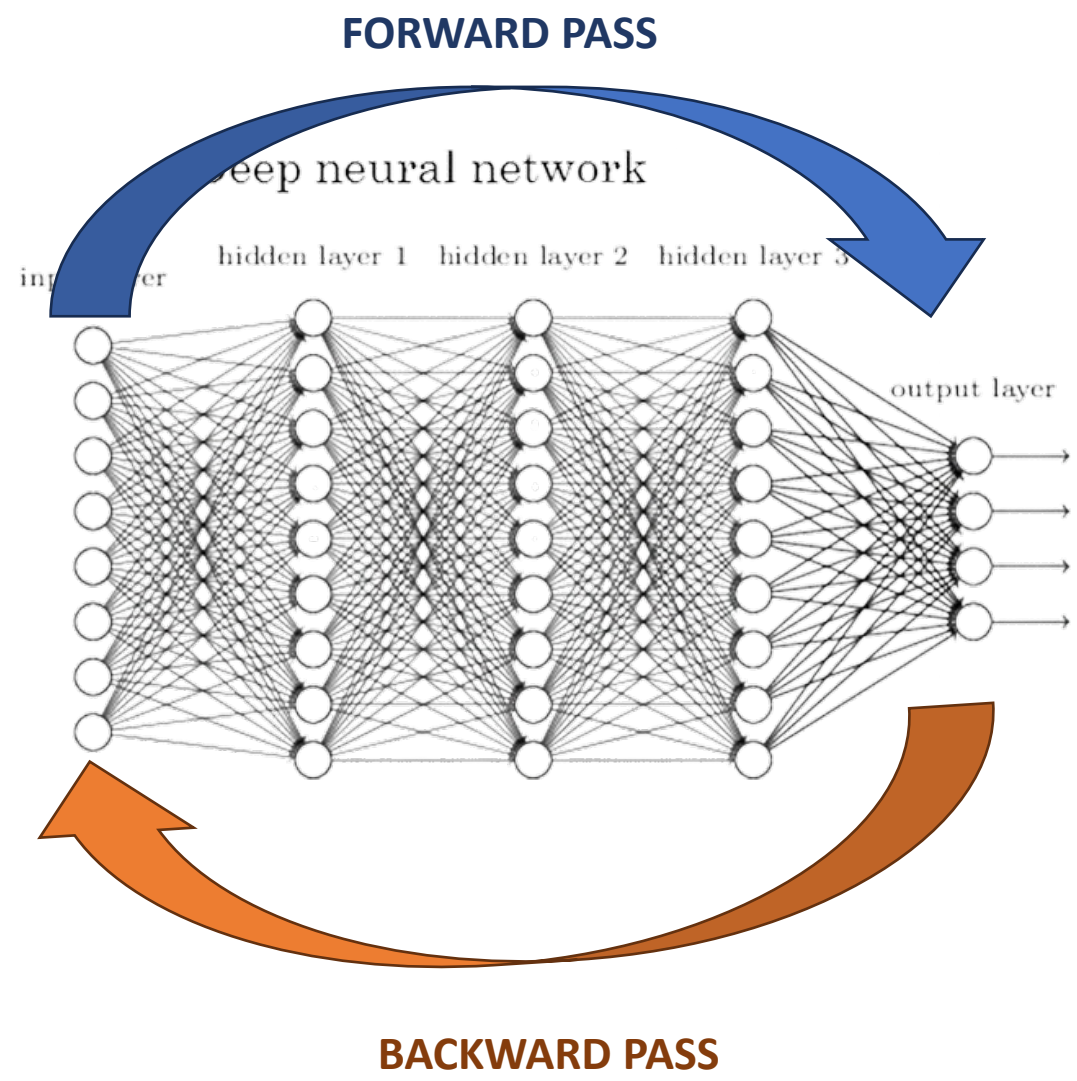
# Entrenamiento

- **Gradient Based Optimization:**
  - Superficie de error compleja
  - Múltiples mínimos locales
  - No hay garantía de llegar al mínimo global



# Entrenamiento

- **BACKPROPAGATION:**
  - Método para implementar el aprendizaje basado en gradiente (GRADIENT BASED)
- Red Neuronal: Función (todo lo compleja que podamos imaginar) de los pesos  $w$  y de la entrada  $x$ :  $o_n = f(w, x)$
- $J(X, Y, \theta)$ : Discrepancia entre la salida obtenida y la esperada (etiquetas) para cada ejemplo.
- **Idea Básica:** Calcular de forma eficiente todas las **derivadas parciales** de la función de coste,  $J(X, Y, \theta)$ , respecto de todos los parámetros ajustables de la función, sus pesos  $w$ , para un valor dado de la entrada  $x$



$$J = \sum_n J_n(o_n, y_n) \qquad \frac{\partial J}{\partial w_i} = \sum_n \frac{\partial J_n}{\partial o_n} \frac{\partial o_n}{\partial w_i}$$



# Entrenamiento

---

- Los pesos se ajustan siguiendo la dirección contraria al gradiente

$$\mathbf{w}_i^{(j+1)} = \mathbf{w}_i^{(j)} - \eta \frac{\partial J}{\partial \mathbf{w}_i} = \mathbf{w}_i^{(j)} - \eta \sum_n \frac{\partial J_n}{\partial \mathbf{o}_n} \frac{\partial \mathbf{o}_n}{\partial \mathbf{w}_i}$$

- Este proceso se realiza de forma iterativa (como en los algoritmos adaptativos) utilizando los datos de entrenamiento
- $\eta$  es lo que se conoce como tasa de aprendizaje o learning rate :
  - Valores grandes: aprendizaje más rápido y menos preciso.
  - Valores pequeños: aprendizaje más lento, pero más preciso

# Entrenamiento

- Es necesario conocer el gradiente para todas y cada una de las capas
- Esta información debe fluir desde la capa de salida hacia la capa de entrada.
- Obtener la expresión analítica del gradiente es posible (sencillo) pero computacionalmente caro
  - Backpropagation hace esto computacionalmente eficiente, aplicando la regla de la cadena (chain rule)

Si  $f$  y  $g$  son funciones reales de variable real:  $y = g(x)$      $z = f(y) = f(g(x))$      $\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$

Si  $\underline{x} \in \mathbb{R}^m$ ;  $\underline{y} \in \mathbb{R}^n$ ;  $\underline{g}$  mapea de  $\mathbb{R}^m$  a  $\mathbb{R}^n$ , y  $f$  mapea de  $\mathbb{R}^n$  a  $\mathbb{R}$

$$\underline{y} = \underline{g}(\underline{x}) \quad z = f(\underline{y}) = f(\underline{g}(\underline{x}))$$

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

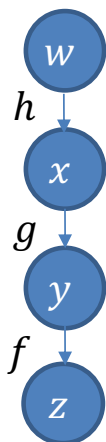


# Computational Graphs

- ¿Como llevar la información del gradiente a todos y cada uno de los pesos de forma eficiente?
  - Se crean grafos de las expresiones involucradas en el cálculo de las salidas
  - Estos grafos indican qué depende de qué y cómo

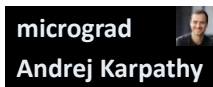
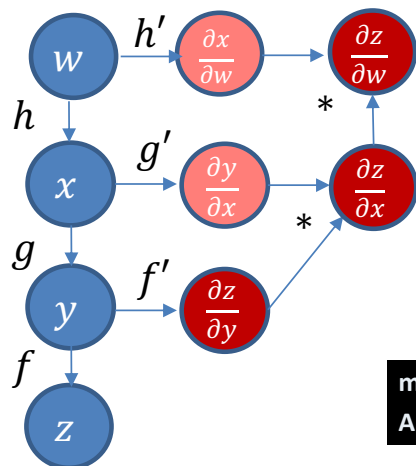
## Forward Pass

$$z = f(y) = f(g(x)) = f(g(h(w)))$$



## Backward Pass

$$\frac{\partial z}{\partial w} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w}$$

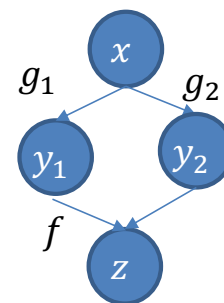


## Forward Pass

$$z = f(y_1, y_2);$$

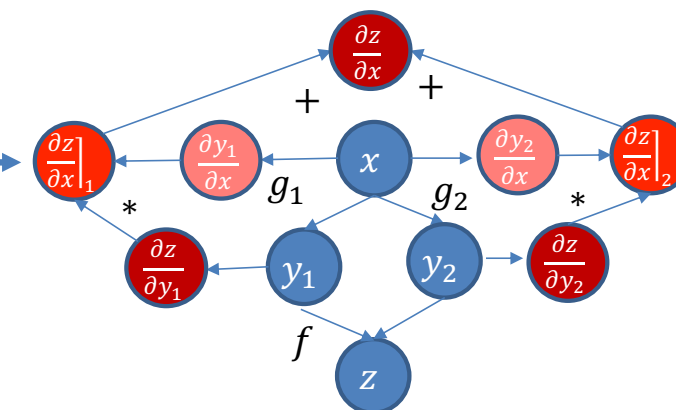
$$y_1 = g_1(x)$$

$$y_2 = g_2(x)$$

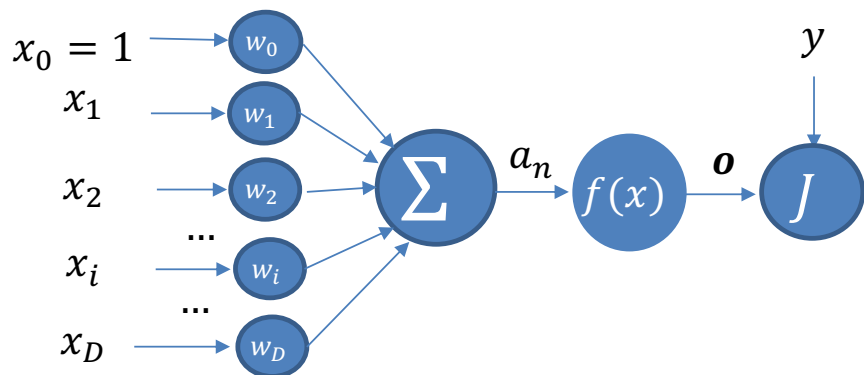


## Backward Pass

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial x} + \frac{\partial z}{\partial y_2} \frac{\partial y_2}{\partial x}$$



# Perceptrón con función de coste MSE:



## Forward Pass

$$a_n = \sum_{i=0}^D w_i x_i = \mathbf{w}^T \mathbf{x}_n$$

$$o_n = f(a_n)$$

$$e_n = y_n - o_n$$

$$J = \frac{1}{2} (e_n)^2$$

## Backward Pass

$$\nabla_{w_i}(J) = \frac{\partial J}{\partial w_i} = \frac{\partial J}{\partial o} \frac{\partial o}{\partial a} \frac{\partial a}{\partial w_i}$$

$$\frac{\partial a}{\partial w_i} = x_i \quad i = 1, 2, \dots, D$$

$$\frac{\partial a}{\partial w_0} = x_0 = 1$$

$$\frac{\partial o}{\partial a} = f'(a) = f'(\mathbf{w}^T \mathbf{x}_n)$$

$$\frac{\partial J}{\partial o} = -(y - o) = -e$$

$$\frac{\partial J}{\partial w_i} = -e \cdot f'(\mathbf{w}^T \mathbf{x}_n) \cdot x_i$$



$$w_i^{(j+1)} = w_i^{(j)} - \eta \frac{\partial J}{\partial w_i} = w_i^{(j)} + \eta \cdot e \cdot f'(\mathbf{w}^T \mathbf{x}_n) \cdot x_i$$

# Stochastic Gradient Descent - SGD

- El gradiente no se calcula haciendo uso de toda la base de datos de entrenamiento sino por bloques y además en orden aleatorio en cada iteración
- SGD converge a mínimos locales. Como la función de coste ya no es tan simple como en el caso lineal, el mínimo local no tiene por qué ser global

- Gradiente Instantáneo:

- Actualización tras cada ejemplo de entrenamiento  $w_i(n+1) = w_i(n) - \eta \frac{\partial J(x_n, y_n, \mathbf{w}^{(n)})}{\partial w_i}$

- Gradiente por lotes, (Batch Gradient):

- Actualización tras cada bloque k  $w_i(k+1) = w_i(k) - \eta \frac{\partial \sum_{n \in k} J(x_n, y_n, \mathbf{w}^{(k)})}{\partial w_i}$

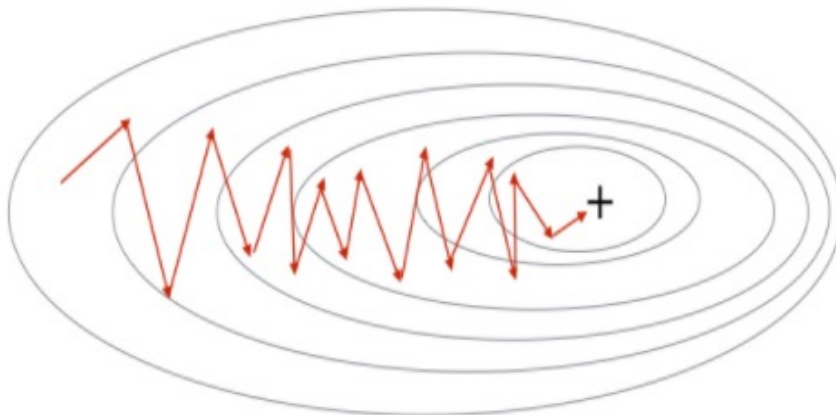
- Gradiente Global:

- Cálculo para toda la base de datos de entrenamiento.
- Actualización tras cada iteración (epoch) de toda la base de datos.
- Alto coste computacional, robusto y baja velocidad de convergencia  $w_i^{(j+1)} = w_i^{(j)} - \alpha \frac{\partial \sum_n J(x_n, y_n, \mathbf{w}^{(j)})}{\partial w_i}$

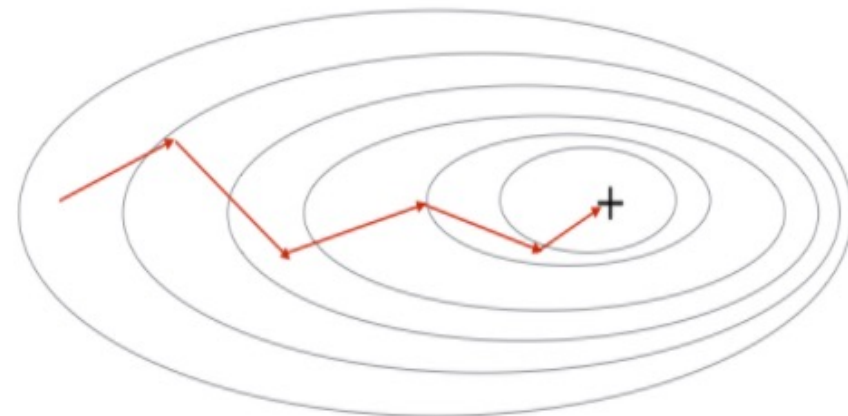
# Backpropagation: Implicaciones y orientaciones prácticas

- **Mini-Batch processing:**
  - Mejor resultado, uso de “pequeños” lotes: **mini-batches**
  - Se puede optimizar mucho con el uso de GPU
  - Es conveniente aleatorizar los datos de entrenamiento
  - Uso de ejemplos más difíciles

Stochastic Gradient Descent



Mini-Batch Gradient Descent



# Implicaciones y orientaciones prácticas

- **Normalización:**

- Acondicionamiento de las entradas (y valores intermedios)
  - Conveniente para un buen funcionamiento
  - Eliminar media y aproximar varianza de todas las dimensiones
  - Decorrelación de las diferentes variables en juego

- Z-NORM

$$\underline{x'} = \frac{\underline{x} - \underline{\mu}_x}{\sigma_x}$$

MIN-MAX

$$\underline{x'} = \frac{\underline{x} - \underline{x}_{min}}{\underline{x}_{max} - \underline{x}_{min}}$$

# Implicaciones y orientaciones prácticas

---

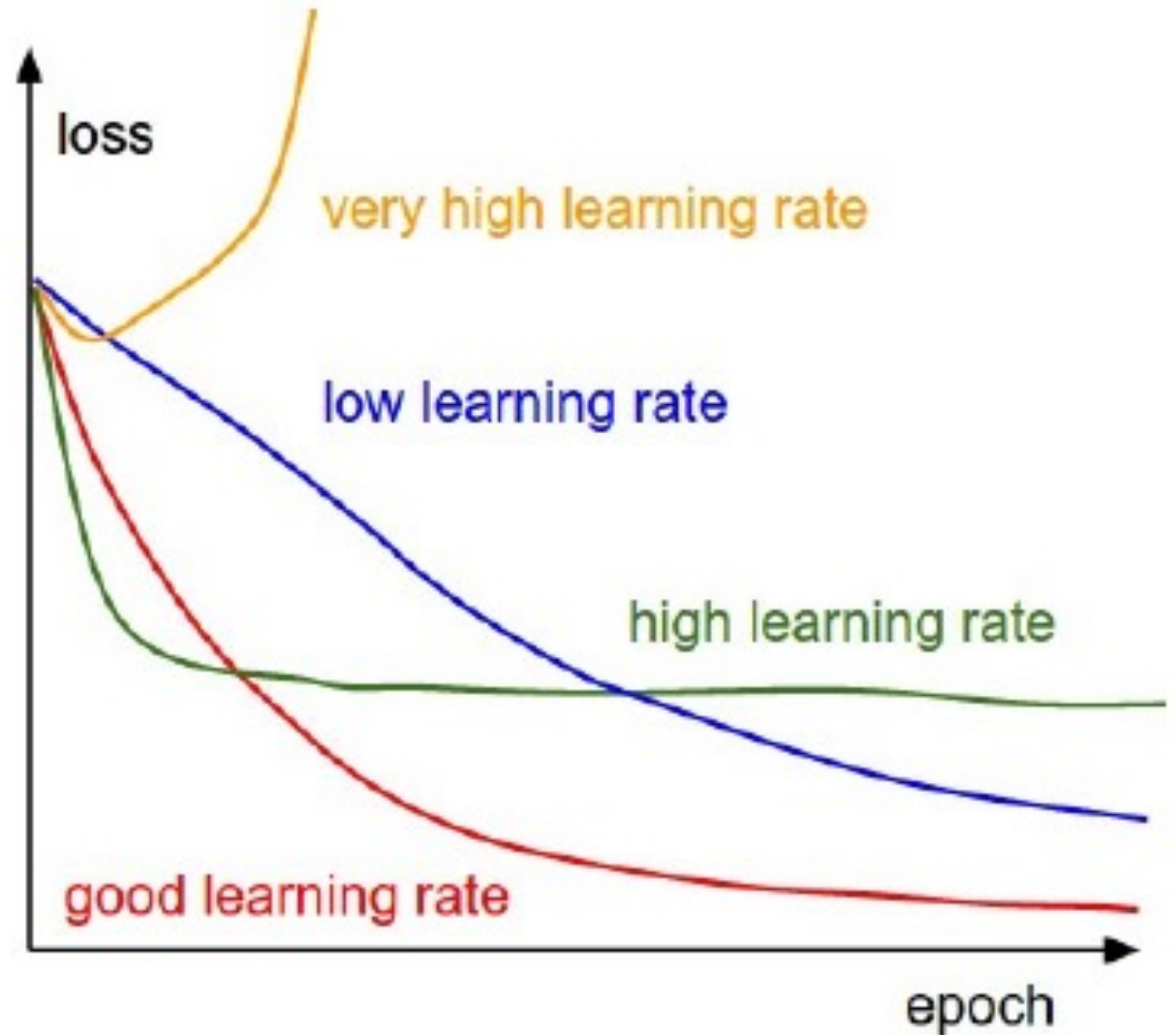
- **INICIALIZACIÓN:**
  - Aleatorizar el valor inicial de los pesos
    - Hay que **“romper la simetría”**
    - Distintas unidades conectadas a las mismas entradas inicializadas con el mismo valor tenderán a aprender lo mismo.
    - Distribución uniforme con media 0 y desviación estándar dependiente de las conexiones que llegan al nodo.
      - pequeñas variaciones en los pesos pueden tener alta repercusión si hay muchas conexiones (fan-in alto)
  - Convergencia vs Generalización:
    - Buena convergencia, llegada a un punto estable rápido
    - Buena generalización, buen funcionamiento para datos no vistos durante el entrenamiento y con propiedades ligeramente distintas
  - Algunas inicializaciones son buenas para una rápida convergencia, pero malas para la generalización y otras al revés.



# Implicaciones y orientaciones prácticas

- **Learning Rate:**

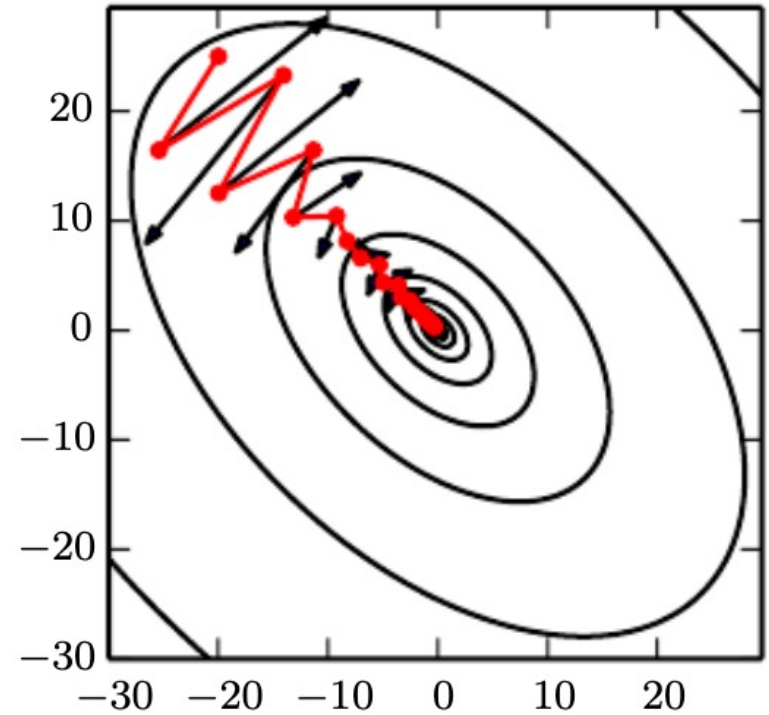
- Valor fijo
- Adaptable
- Ajustado a cada conexión



# Estrategias de Aceleración de Convergencia

- **Momentos:**

- El gradiente modifica la “velocidad” de actualización de los pesos
- Se le aporta “inercia” al proceso
- Se acumula una media móvil de gradientes pasados para seguir moviéndose en esa dirección



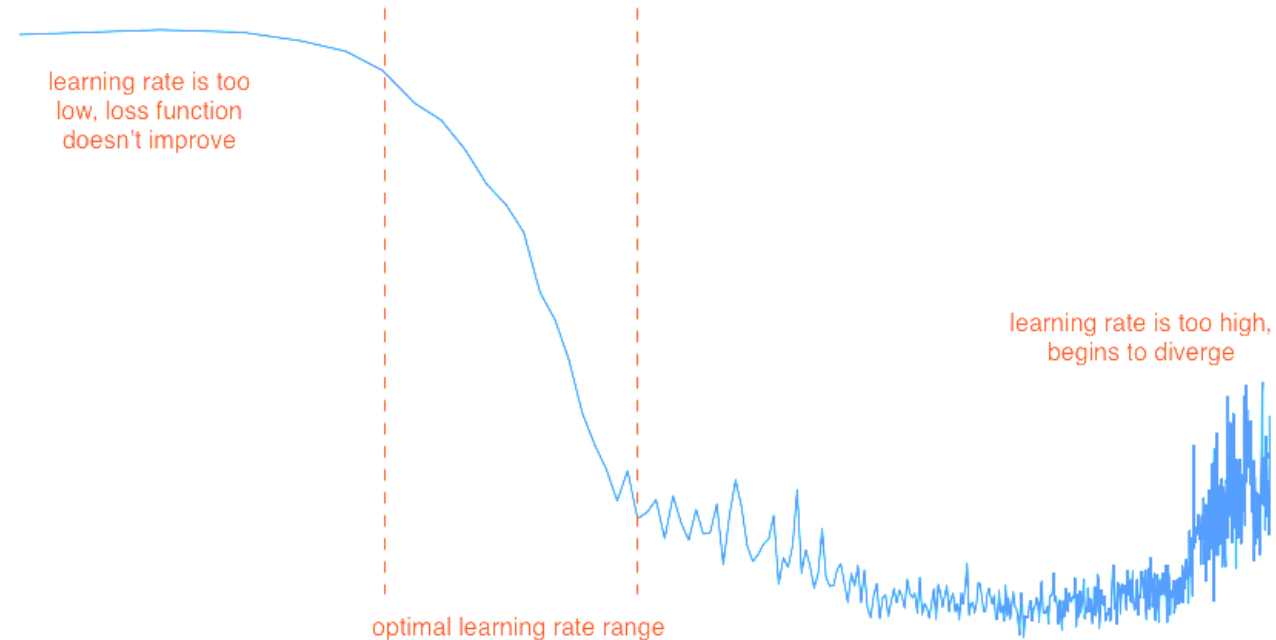
$$v_i(k) = \beta \cdot v_i(k-1) - \eta \frac{\partial \sum_{n \in k} \mathcal{J}(x_n, y_n, w^{(k)})}{\partial w_i}$$

$$w_i(k+1) = w_i(k) + v_i(k)$$

# Estrategias de Aceleración de Convergencia

- **Adaptive Learning Rate:**

- Se reduce el learning rate cuando la reducción de la función de coste se ralentiza (plateau)
- Se suele usar SGD con momentum y decaying learning rate (hasta un valor mínimo del learning rate, con caída exponencial, o dividiéndolo por un factor entre 2 y 10 cuando se detecta la meseta).



# Estrategias de Aceleración de Convergencia

- **Adaptive Learning Rate:**

- **AdaGrad**

- Se acumulan los gradientes al cuadrado

$$\mathbf{r}(k) = \mathbf{r}(k - 1) + \left( \frac{\partial \sum_{n \in k} \mathcal{J}(x_n, y_n, \mathbf{w}^{(k)})}{\partial \mathbf{w}_i} \right)^2$$

- El learning rate se divide por la raíz cuadrada del acumulado de los gradientes al cuadrado

$$\eta(k) = \frac{\epsilon}{\delta + \sqrt{\mathbf{r}(k)}}$$

# Estrategias de Aceleración de Convergencia

- **Adaptive Learning Rate:**

- **RMSProp**

- Modificación sobre AdaGrad
- Cambia la acumulación de gradientes al cuadrado por una media móvil (exponencial)

$$\mathbf{r}(k) = \rho \cdot \mathbf{r}(k - 1) + (1 - \rho) \left( \frac{\partial \sum_{n \in k} \mathcal{J}(\mathbf{x}_n, \mathbf{y}_n, \mathbf{w}^{(k)})}{\partial \mathbf{w}_i} \right)^2$$

- El learning rate se divide por la raíz cuadrada del acumulado de los gradientes al cuadrado

$$\eta(k) = \frac{\epsilon}{\delta + \sqrt{\mathbf{r}(k)}}$$

# Estrategias de Aceleración de Convergencia

- **Adaptive Learning Rate:**

- **Adam**

- Añade momento al RMSProp (Adam: “adaptive moments”)

$$s_i(k) = \rho_1 \cdot s_i(k-1) + (1 - \rho_1) \frac{\partial \sum_{n \in k} \mathcal{J}(x_n, y_n, \mathbf{w}^{(k)})}{\partial w_i}$$

- Acumulación de gradientes al cuadrado con una media móvil (exponencial) para adaptar Learning rate

$$r(k) = \rho_2 \cdot r(k-1) + (1 - \rho_2) \left( \frac{\partial \sum_{n \in k} \mathcal{J}(x_n, y_n, \mathbf{w}^{(k)})}{\partial w_i} \right)^2 \quad \eta(k) = \frac{\epsilon}{\delta + \sqrt{r(k)}}$$

- Combinados

$$w_i(k+1) = w_i(k) - \eta(k) s_i(k)$$



# Estrategias de Aceleración de Convergencia

---

- Momentos
- Adaptive Learning Rate
  - AdaGrad
  - RMSprop
  - Adam
  - ...

- ¿Cuál usar?

No hay consenso claro y depende del problema y de los datos.

# Implicaciones y orientaciones prácticas

- **Batch Normalization:**
  - Estrategia para reparametrizar redes profundas
  - Los gradientes y actualizaciones en una capa dependen en gran medida del resto de las capas lo que hace difícil el ajuste de hiperparámetros como el learning rate.
  - Si  $H$ , es un minibatch de activaciones de una capa, estas se normalizan para que tengan media nula y varianza unidad

$$H' = \frac{H - \mu}{\sigma} \quad \mu = \frac{1}{m} \sum_i H_i \quad \sigma = \sqrt{\delta + \frac{1}{m} \sum_i (H - \mu)_i^2}$$

- En test  $\mu$  y  $\sigma$  pueden tomar valores aprendidos durante el entrenamiento



# Implicaciones y orientaciones prácticas

- **Técnicas de Generalización y Regularización:**

- **Generalización:**

- Si la red está sobredimensionada (demasiados parámetros) se corre el riesgo de aprender las peculiaridades de los datos de entrenamiento:

- Sobreajuste u overfitting.

- Para evitar el sobreajuste:

- Contar con más datos:

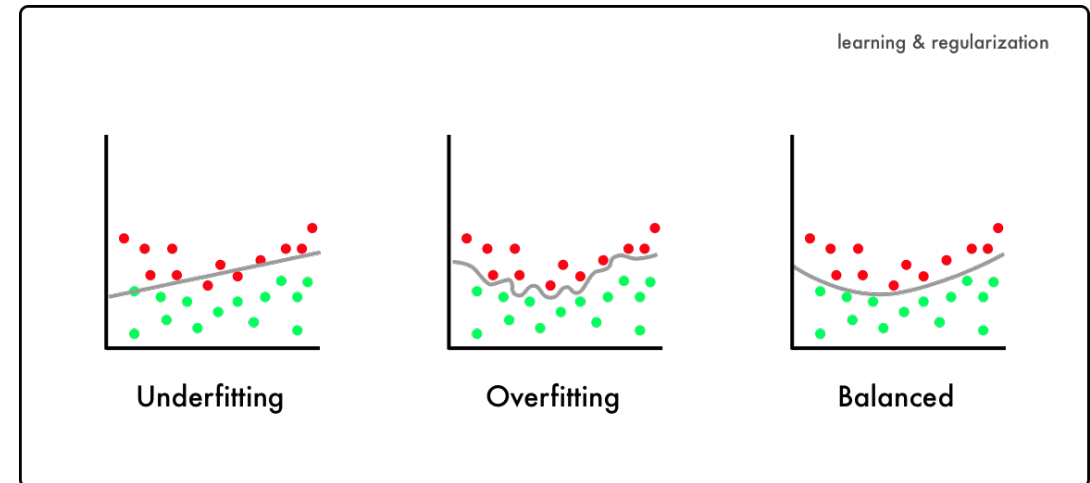
- Aumentar los dataset de train

- **Data augmentation**

- Limitar la capacidad de la red para no exceder la complejidad del problema que se aborda:

- Limitar el número de capas

- **Weight Sharing:** Varias unidades comparten sus pesos



# Implicaciones y orientaciones prácticas

- **Técnicas de Generalización y Regularización :**
  - Estrategias:
    - **Early stopping:**
      - El entrenamiento se detiene antes de que la red llegue a overfitting
    - **Weight decay:**
      - Se intenta que los pesos no alcancen valores muy grandes (norma L2 ó norma L1)
    - **Ruido:**
      - Se añade un pequeño ruido en algunos puntos de la red durante el entrenamiento
    - **Dropout:**
      - Se desactivan algunas unidades durante el entrenamiento para que el resto se adapte en su ausencia.