

**2020-2021**

## **Aprendizaje Automático**

# **5. Redes Neuronales Multicapa**



Francisco Casacuberta Nolla  
(fcn@dsic.upv.es)

Enrique Vidal Ruiz  
(evidal@dsic.upv.es)

Departament de Sistemes Informàtics i Computació (DSIC)

Universitat Politècnica de València (UPV)

# Index

- 1 Redes neuronales multicapa ▷ 2
- 2 Algoritmo de retropropagación del error (BackProp) ▷ 18
- 3 Aspectos de uso y propiedades del BackProp ▷ 33
- 4 Variantes de BackProp ▷ 41
- 5 Redes neuronales radiales ▷ 47
- 6 Aplicaciones ▷ 52
- 7 Notación ▷ 54

# Index

- 1 *Redes neuronales multicapa* ▷ 2
- 2 Algoritmo de retropropagación del error (BackProp) ▷ 18
- 3 Aspectos de uso y propiedades del BackProp ▷ 33
- 4 Variantes de BackProp ▷ 41
- 5 Redes neuronales radiales ▷ 47
- 6 Aplicaciones ▷ 52
- 7 Notación ▷ 54

# Modelo conexionista

- Un conjunto de procesadores elementales densamente interconectados.
- Nombres alternativos:
  - Modelo conexionista.
  - Red neuronal artificial.
  - Procesado distribuido y paralelo.
- Perceptrón multicapa:
  - Modelo conexionista simple.
  - Fronteras de decisión complejas.
  - Optimización no convexa.
  - Entrenamiento de los pesos mediante descenso por gradiente: algoritmo de **retropropagación del error**.

## Introducción: historia (I)

- 1943: McCulloch y Pitt introducen un modelo matemático simple de “neurona”.
- 1949: Hebb propone una regla que modela el aprendizaje en las neuronas. Rochester realiza una simulación en un computador IBM en 1950.
- 1957: *Rosenblatt* introduce *el Perceptrón* como un dispositivo hardware con capacidad de autoaprendizaje y Widrow y Hoff proponen el *Adaline* para la cancelación de ecos en redes telefónicas.
- 1969: *Minsky y Papert* demuestran que un perceptrón solo puede implementar funciones discriminantes lineales y que esta limitación no se puede superar mediante multiples perceptrones organizados en cascada: Para ello sería necesario introducir funciones no-lineales.
- 1970-1975: Diversos autores tratan de desarrollar algoritmos de descenso por gradiente adecuados para multiples perceptrones en cascada con funciones no-lineales. El cálculo de derivadas parciales se muestra esquivo.

## Introducción: historia (II)

- 1986: *Rumelhart, Hinton y Williams* popularizan la técnica de *retropropagación del error*. Se basa en el uso de cierto tipo de funciones no lineales, llamadas “funciones de activación”, con las que se simplifica el cálculo de derivadas parciales necesarias para descenso por gradiente. Al parecer, técnicas similares habían sido ya propuestas por *Linnainmaa* en 1970 y *Werbos* en 1974.
- 1996: *Bishop, Rippley, Ney*, entre otros, dan una interpretación probabilística a las redes neuronales y al algoritmo de retropropagación del error.
- 2006: *Hinton* publica en *Science* un artículo que inaugura una nueva tendencia denominada “*redes profundas*”. Posteriormente, en 2015 *LeCun, Bengio y Hinton* publican un artículo sobre estas técnicas en *Nature*. Se desarrollan diversas técnicas que mejoran el uso del algoritmo de retropropagación en redes profundas y en redes recurrentes. Aplicadas con gran éxito en múltiples problemas.

# Funciones discriminantes lineales y función de activación

- FUNCIONES DISCRIMINANTES LINEALES (FDL)

$$\phi : \mathbb{R}^d \rightarrow \mathbb{R} : \phi(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\theta}^t \mathbf{x} = \sum_{i=0}^d \theta_i x_i$$

Notación en coordenadas homogéneas:

- $\mathbf{x} \in \mathbb{R}^{d+1}$ ,  $\mathbf{x} = x_0, x_1, \dots, x_d$ ,  $x_0 \stackrel{\text{def}}{=} 1$
- $\boldsymbol{\theta} \in \mathbb{R}^D$ ,  $\boldsymbol{\theta} = \theta_0, \theta_1, \dots, \theta_d$   $D \stackrel{\text{def}}{=} d + 1$

La componente 0 del *vector de pesos* es el *umbral*,  $\theta_0 \in \mathbb{R}$

# Funciones discriminantes lineales y función de activación

- FUNCIONES DISCRIMINANTES LINEALES (FDL)

$$\phi : \mathbb{R}^d \rightarrow \mathbb{R} : \phi(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\theta}^t \mathbf{x} = \sum_{i=0}^d \theta_i x_i$$

Notación en coordenadas homogéneas:

$$\begin{aligned} - \mathbf{x} &\in \mathbb{R}^{d+1}, & \mathbf{x} &= x_0, x_1, \dots, x_d, & x_0 &\stackrel{\text{def}}{=} 1 \\ - \boldsymbol{\theta} &\in \mathbb{R}^D, & \boldsymbol{\theta} &= \theta_0, \theta_1, \dots, \theta_d & D &\stackrel{\text{def}}{=} d+1 \end{aligned}$$

La componente 0 del *vector de pesos* es el *umbral*,  $\theta_0 \in \mathbb{R}$

- FUNCIONES DISCRIMINANTES LINEALES CON ACTIVACIÓN (FDLA)

$$g \circ \phi : \mathbb{R}^d \rightarrow \mathbb{R} : g \circ \phi(\mathbf{x}; \boldsymbol{\theta}) = g(\boldsymbol{\theta}^t \mathbf{x})$$

$g : \mathbb{R} \rightarrow \mathbb{R}$  es una *función de activación*<sup>†</sup>

<sup>†</sup> también denominada *función logística* y a la FDLA *función discriminante lineal logística*.



# Funciones de activación y sus derivadas

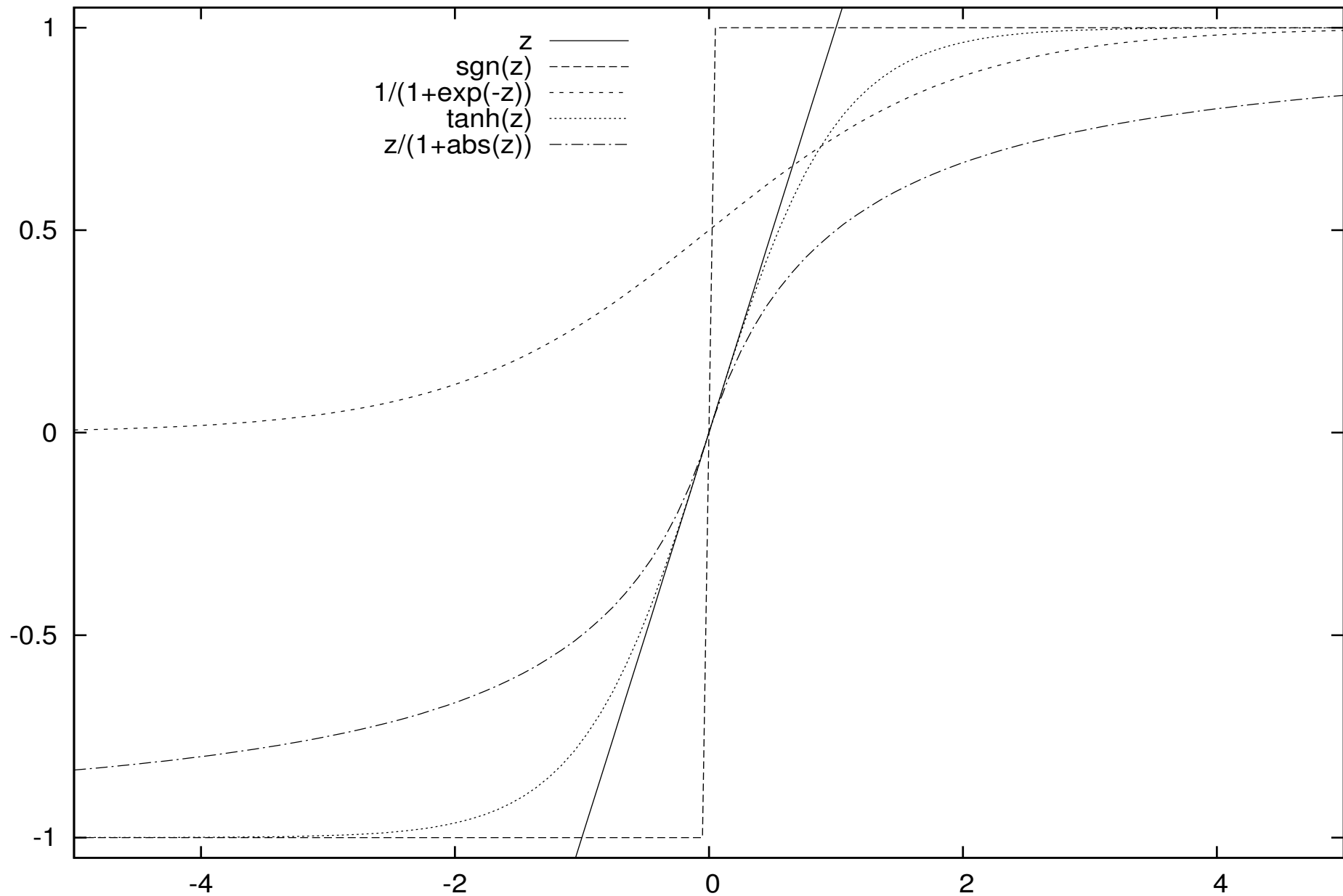
Sea  $g : \mathbb{R} \rightarrow \mathbb{R}$  y  $z \in \mathbb{R}$ :

- **LINEAL**:  $g_L(z) = z \Rightarrow g'_L(z) = \frac{d g_L}{d z} = 1$
- **RELU** (rectified linear unit):  $g_U(z) = \max(0, z) \Rightarrow g'_U(z) = \begin{cases} 0 & \text{si } z < 0 \\ 1 & \text{si } z > 0 \\ \text{no definida} & \text{si } z = 0 \end{cases}$
- **ESCALÓN**<sup>†</sup>:  $g_E(z) = \text{sgn}(z) \stackrel{\text{def}}{=} \begin{cases} +1 & \text{si } z > 0 \\ -1 & \text{si } z < 0 \end{cases} \Rightarrow g'_E(z) = \begin{cases} \text{no definida} & \text{si } z = 0 \\ 0 & \text{si } z \neq 0 \end{cases}$
- **SIGMOID**:  $g_S(z) = \frac{1}{1 + \exp(-z)} \Rightarrow g'_S(z) = g_S(z) (1 - g_S(z))$
- **TANGENTE HIPERBÓLICA**:  $g_T(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)} \Rightarrow g'_T(z) = 1 - (g_T(z))^2$
- **RÁPIDA**:  $g_F(z) = \frac{z}{1 + |z|} \Rightarrow g'_F(z) = \frac{1}{(1 + |z|)^2} = \left( \frac{g_F(z)}{z} \right)^2$
- **SOFTMAX**: Para  $z_1, \dots, z_n \in \mathbb{R}$ ,  $g_M(z_i) = \frac{\exp(z_i)}{\sum_j \exp(z_j)} \Rightarrow g'_M(z_i) = g_M(z_i) (1 - g_M(z_i))$

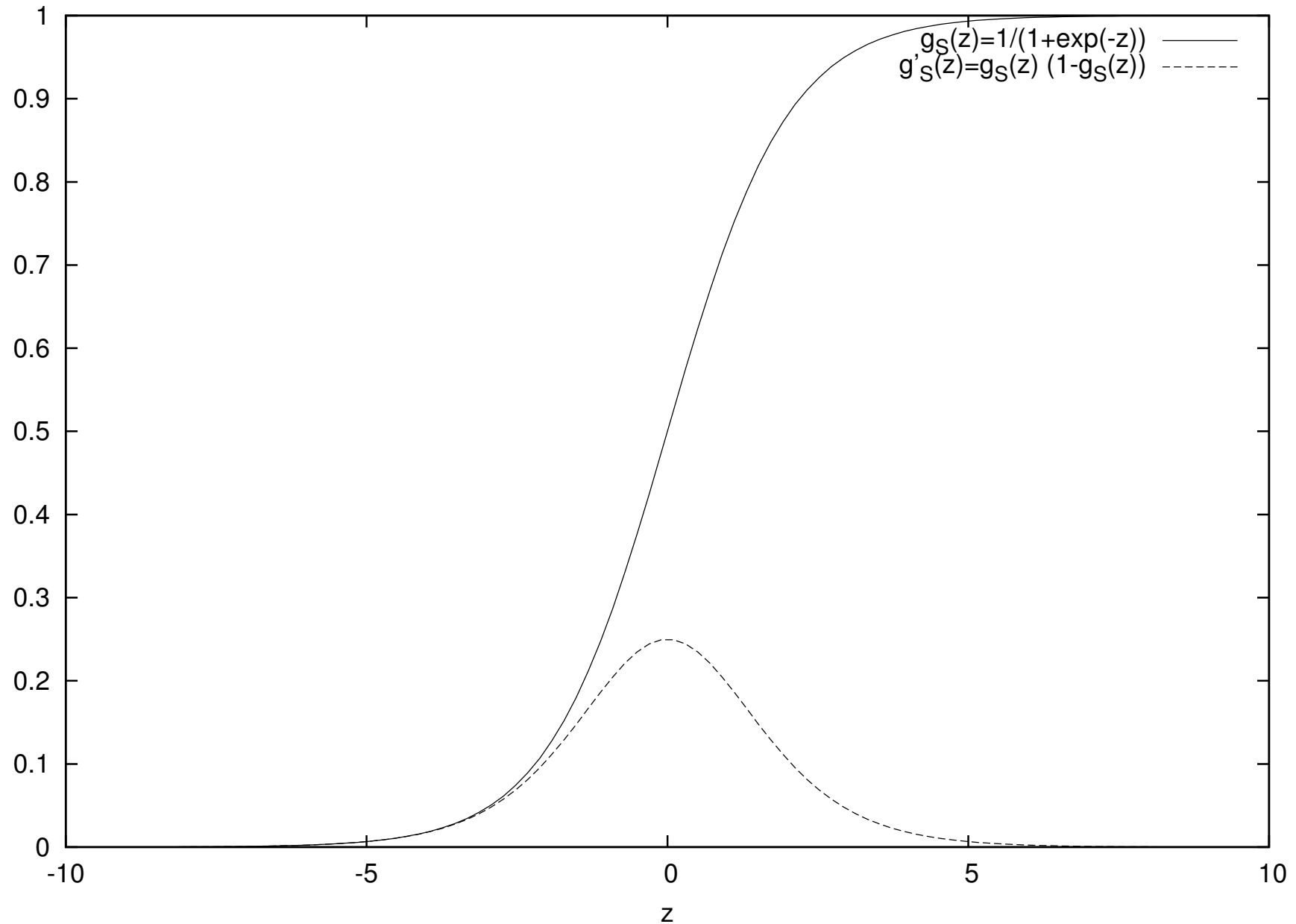
Ejercicios: Demostrar  $g'_S(z) = g_S(z) (1 - g_S(z))$ ,  $g'_T(z) = 2g_S(2z) - 1 \quad \forall z \in \mathbb{R}$

<sup>†</sup>  $\text{sgn}$  es la *función signo*

# Funciones de activación

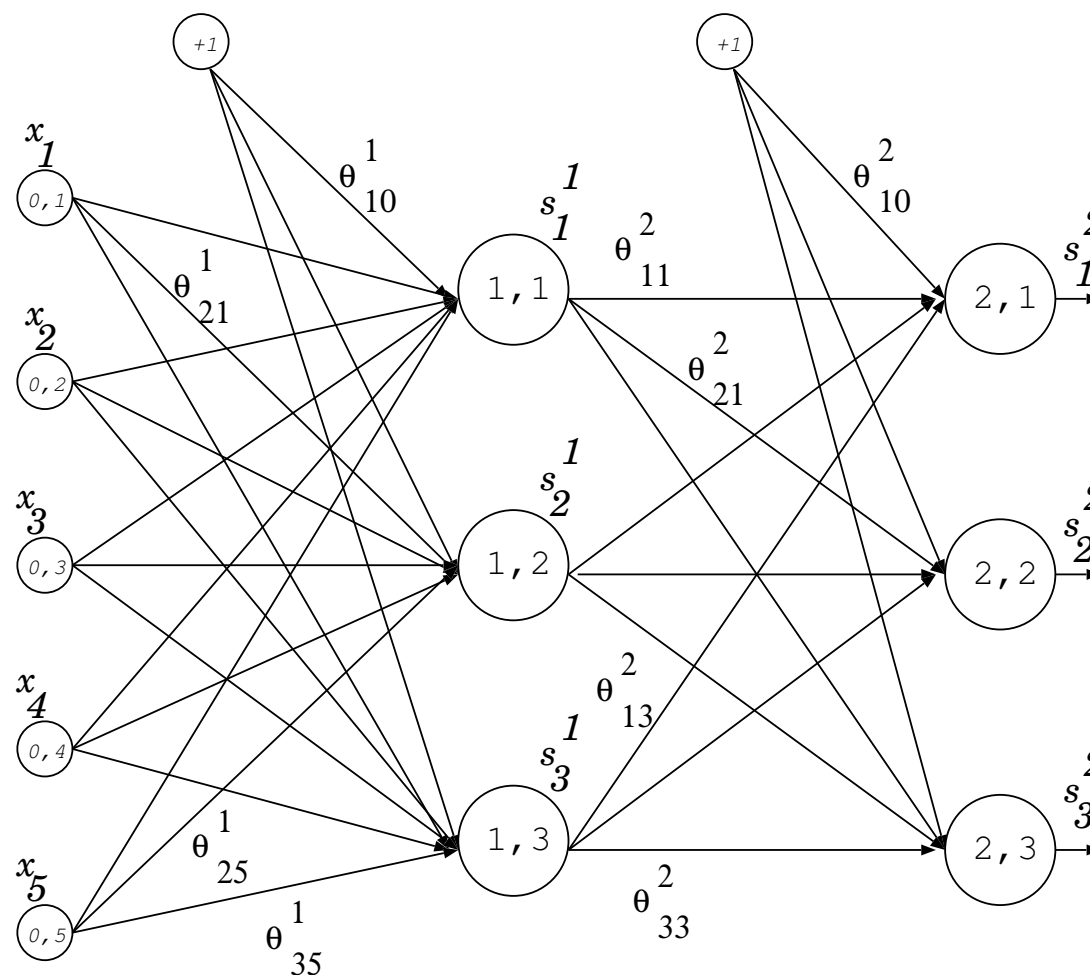


# Derivada de la función de activación sigmoid



# Un perceptrón de dos capas: ejemplo y notación

Topología



Dinámica

Capa de entrada

$$1 \leq i \leq M_0 \equiv d = 5$$

$$x_i \in \mathbb{R}$$

Capa oculta

$$1 \leq i \leq M_1 = 3$$

$$s_i^1(\mathbf{x}; \Theta) = g\left(\sum_{j=0}^{M_0} \theta_{ij}^1 x_j\right)$$

Capa de salida

$$1 \leq i \leq M_2 = 3$$

$$s_i^2(\mathbf{x}; \Theta) = g\left(\sum_{j=0}^{M_1} \theta_{ij}^2 s_j^1(\mathbf{x})\right)$$

## Perceptrón de dos capas

- *Un perceptrón de dos capas* consiste en una combinación de FDLA agrupadas en 2 capas de tallas  $M_1$  (capa *oculta*) y  $M_2$  (capa de salida), más una capa de entradas de talla  $M_0 = d$  (por simplicidad no se contabilizan los umbrales):

$$s_i^2(\mathbf{x}; \Theta) = g\left(\sum_{j=0}^{M_1} \theta_{ij}^2 s_j^1(\mathbf{x}; \Theta)\right) = g\left(\sum_{j=0}^{M_1} \theta_{ij}^2 g\left(\sum_{j'=0}^{M_0} \theta_{jj'}^1 x_{j'}\right)\right) \quad 1 \leq i \leq M_2$$

Los parámetros son  $\Theta = [\theta_{10}^1, \dots, \theta_{M_1 M_0}^1, \theta_{10}^2, \dots, \theta_{M_2 M_1}^2] \in \mathbb{R}^D$

- *Problema*: Dado un conjunto de entrenamiento  $S = \{(\mathbf{x}_1, \mathbf{t}_1), \dots, (\mathbf{x}_N, \mathbf{t}_N)\}$ , con  $\mathbf{x}_n \in \mathbb{R}^{M_0}$ ,  $\mathbf{t}_n \in \mathbb{R}^{M_2}$ , encontrar  $\Theta$  tal que  $s^2(\mathbf{x}_n; \Theta)$  aproxime lo mejor posible a  $\mathbf{t}_n \forall n, 1 \leq n \leq N$ .
- *En clasificación*:  $M_2 \equiv C$  y las etiquetas  $\mathbf{t}_n$  para  $1 \leq n \leq N$  son de la forma

$$1 \leq c \leq C \quad t_{nc} = \begin{cases} 1 & \mathbf{x}_n \text{ es de la clase } c \\ 0 \text{ (o } -1) & \mathbf{x}_n \text{ no es de la clase } c \end{cases}$$

*Simplificaciones de notación*:  $s_i^k(\mathbf{x}; \Theta) \equiv s_i^k(\mathbf{x}) \equiv s_i^k \quad \forall k, i$

# El perceptrón multicapa y las funciones de activación

- Un perceptrón multicapa de dos capas define una función  $\mathbb{R}^{M_0} \rightarrow \mathbb{R}^{M_2}$ :

$$s_i^2(\mathbf{x}) = g\left(\sum_{j=0}^{M_1} \theta_{ij}^2 s_j^1(\mathbf{x})\right) = g\left(\sum_{j=0}^{M_1} \theta_{ij}^2 g\left(\sum_{j'=0}^{M_0} \theta_{jj'}^1 x_{j'}\right)\right) \text{ para } 1 \leq i \leq M_2$$

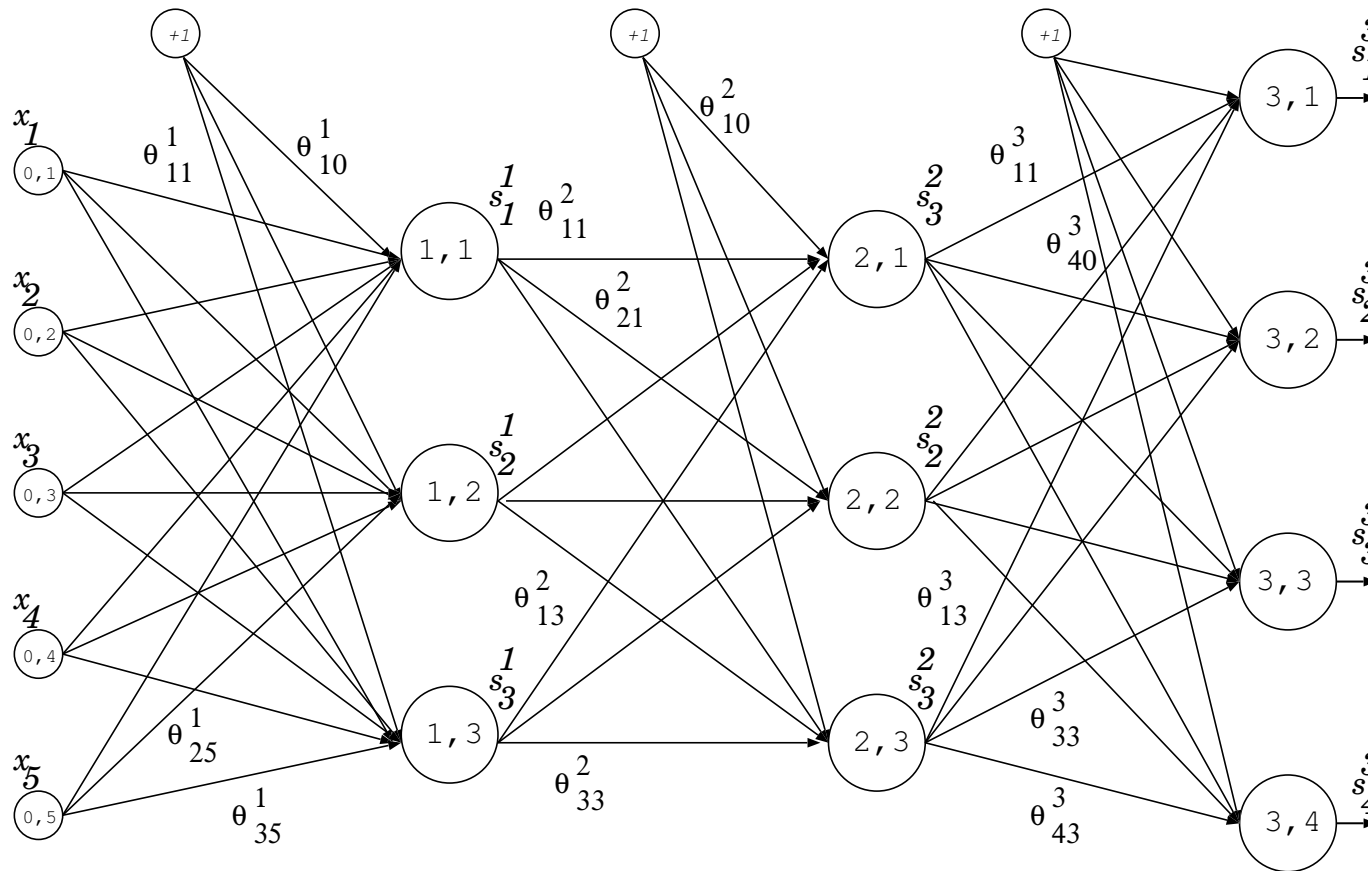
- Si todas las funciones de activación son lineales, un perceptrón multicapa define **UNA FUNCIÓN DISCRIMINANTE LINEAL**,  $\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_{M_2}(\mathbf{x}))^t$ :

$$\phi_i(\mathbf{x}) \equiv s_i^2(\mathbf{x}) = \sum_{j=0}^{M_1} \sum_{j'=0}^{M_0} \theta_{ij}^2 \theta_{jj'}^1 x_{j'} = \sum_{j'=0}^{M_0} \left(\sum_{j=0}^{M_1} \theta_{ij}^2 \theta_{jj'}^1\right) x_{j'} = \sum_{j'=0}^{M_0} \theta_{ij'} x_{j'}$$

- Si al menos una función de activación no es lineal (y, sin pérdida de generalidad, todas las funciones de activación de la capa de salida son lineales) un perceptrón multicapa define **UNA FUNCIÓN DISCRIMINANTE LINEAL GENERALIZADA**:

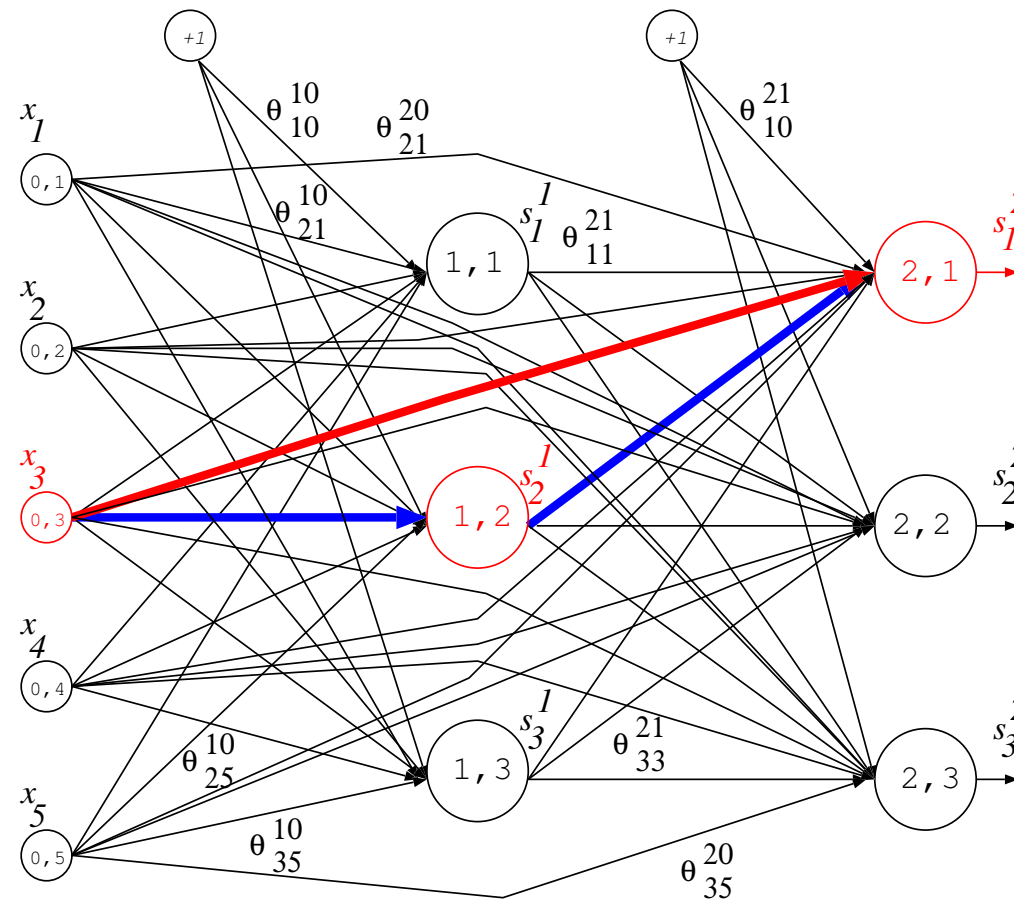
$$\phi_i(\mathbf{x}) \equiv s_i^2(\mathbf{x}) = \sum_{j=0}^{M_1} \theta_{ij}^2 g\left(\sum_{j'=0}^{M_0} \theta_{jj'}^1 x_{j'}\right) = \sum_{j=0}^{M_1} \theta_{ij}^2 \psi_j(\mathbf{x}) \text{ para } 1 \leq i \leq M_2$$

# Un perceptrón de tres capas



- Primera capa oculta:  $s_i^1 = g(\sum_{j=0}^{M_0} \theta_{ij}^1 x_j)$  para  $1 \leq i \leq M_1$
- Segunda capa oculta:  $s_i^2 = g(\sum_{j=0}^{M_1} \theta_{ij}^2 s_j^1)$  para  $1 \leq i \leq M_2$
- Capa de salida:  $s_i^3 = g(\sum_{j=0}^{M_2} \theta_{ij}^3 s_j^2)$  para  $1 \leq i \leq M_3$

# Redes hacia adelante de dos capas



- Primera capa oculta:  $s_i^1 = g(\sum_{j=0}^{M_0} \theta_{ij}^{1,0} x_j)$  para  $1 \leq i \leq M_1$
- Capa de salida:  $s_i^2 = g(\sum_{j=0}^{M_1} \theta_{ij}^{2,1} s_j^1 + \sum_{j=0}^{M_0} \theta_{ij}^{2,0} x_j)$  para  $1 \leq i \leq M_2$

El perceptrón multicapa es un caso particular



# El perceptrón multicapa como regresor

Regresión de  $\mathbb{R}^d$  a  $\mathbb{R}^{d'}$

(p.e. un PM de dos capas con  $d \equiv M_0$  y  $d' \equiv M_2$ )

$$\mathbf{f} : \mathbb{R}^{M_0} \rightarrow \mathbb{R}^{M_2} : f_i(\mathbf{x}) \equiv s_i^2(\mathbf{x}) = \sum_{j=0}^{M_1} \theta_{kj}^2 g\left(\sum_{j'=0}^{M_0} \theta_{jj'}^1 x_{j'}\right) \quad 1 \leq i \leq M_2$$

- Cualquier función se puede aproximar con precisión arbitraria mediante un perceptrón de *una* o más capas ocultas con un número de nodos suficientemente grande
- En general, para alcanzar una precisión dada, el número de nodos necesarios suele ser *mucho menor* si el número de capas ocultas es mayor o igual que *dos*

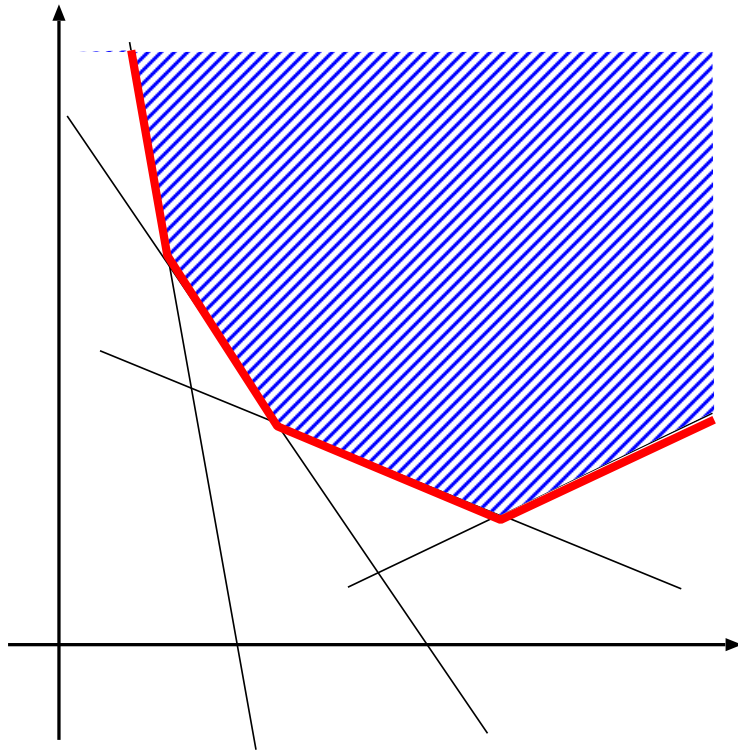
# El perceptrón multicapa como clasificador

**Clasificación** en  $C$  clases de puntos de  $\mathbb{R}^d$  (PM con  $M_0 \equiv d$ ,  $M_2 \equiv C$ )

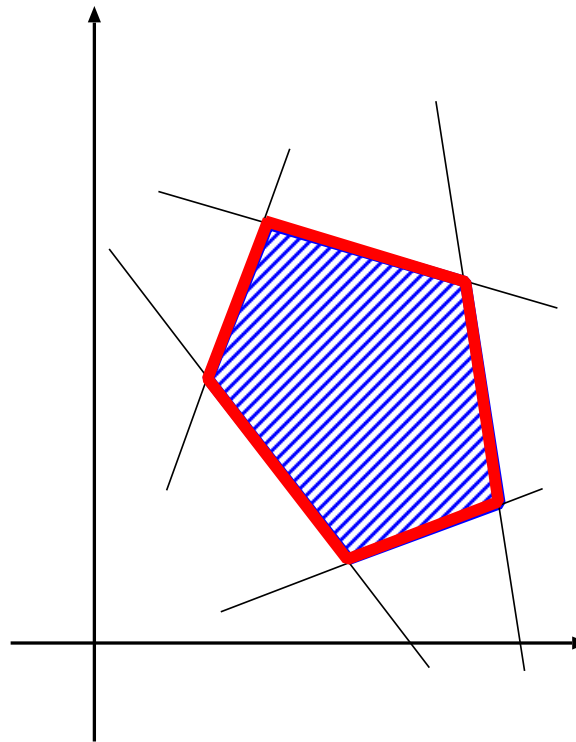
$$f : \mathbb{R}^d \rightarrow \{1, \dots, C\} : f(\mathbf{x}) = \arg \max_{1 \leq c \leq C} \phi_c(\mathbf{x}) = \arg \max_{1 \leq c \leq M_2} s_c^2(\mathbf{x})$$

- Si un conjunto de entrenamiento es linealmente separable existe un perceptrón sin capas ocultas que lo clasifica correctamente.
- Un PM con *una* capa oculta de  $N - 1$  nodos puede clasificar correctamente las muestras de cualquier conjunto de entrenamiento de talla  $N$ . ¿Con qué poder de generalización?
- Cualquier frontera de decisión basada en trozos de hiperplanos puede obtenerse mediante un PM con *una* capa oculta y un número de nodos adecuado [Huang & Lippmann, 1988], [Huang, Chen & Babri, 2000].
- En general, el número de nodos necesarios para aproximar una frontera dada suele ser *mucho menor* si el número de capas ocultas es mayor o igual que *dos*.

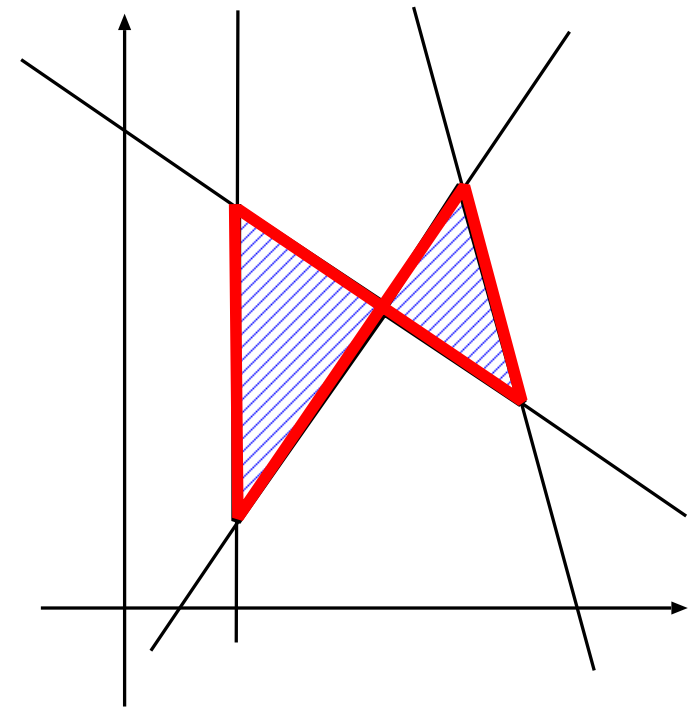
# El perceptrón multicapa como clasificador



Fronteras lineales a intervalos



Regiones cerradas



Regiones no convexas

# Index

- 1 Redes neuronales multicapa ▷ 2
- 2 *Algoritmo de retropropagación del error (BackProp)* ▷ 18
- 3 Aspectos de uso y propiedades del BackProp ▷ 33
- 4 Variantes de BackProp ▷ 41
- 5 Redes neuronales radiales ▷ 47
- 6 Aplicaciones ▷ 52
- 7 Notación ▷ 54

# Aprendizaje de los pesos de un perceptrón multicapa

PROBLEMA (REGRESIÓN): dada la topología de un perceptrón multicapa con  $L$  capas y un conjunto de entrenamiento:

$$S = \{(\mathbf{x}_1, \mathbf{t}_1), \dots, (\mathbf{x}_N, \mathbf{t}_N)\}, \quad \mathbf{x}_n \in \mathbb{R}^{M_0}, \quad \mathbf{t}_n \in \mathbb{R}^{M_L}$$

obtener  $\Theta$  que minimice el error cuadrático medio:

$$q_S(\Theta) = \frac{1}{N} \sum_{n=1}^N q_n(\Theta); \quad q_n(\Theta) = \frac{1}{2} \sum_{i=1}^{M_L} (t_{ni} - s_i^L(\mathbf{x}_n; \Theta))^2$$

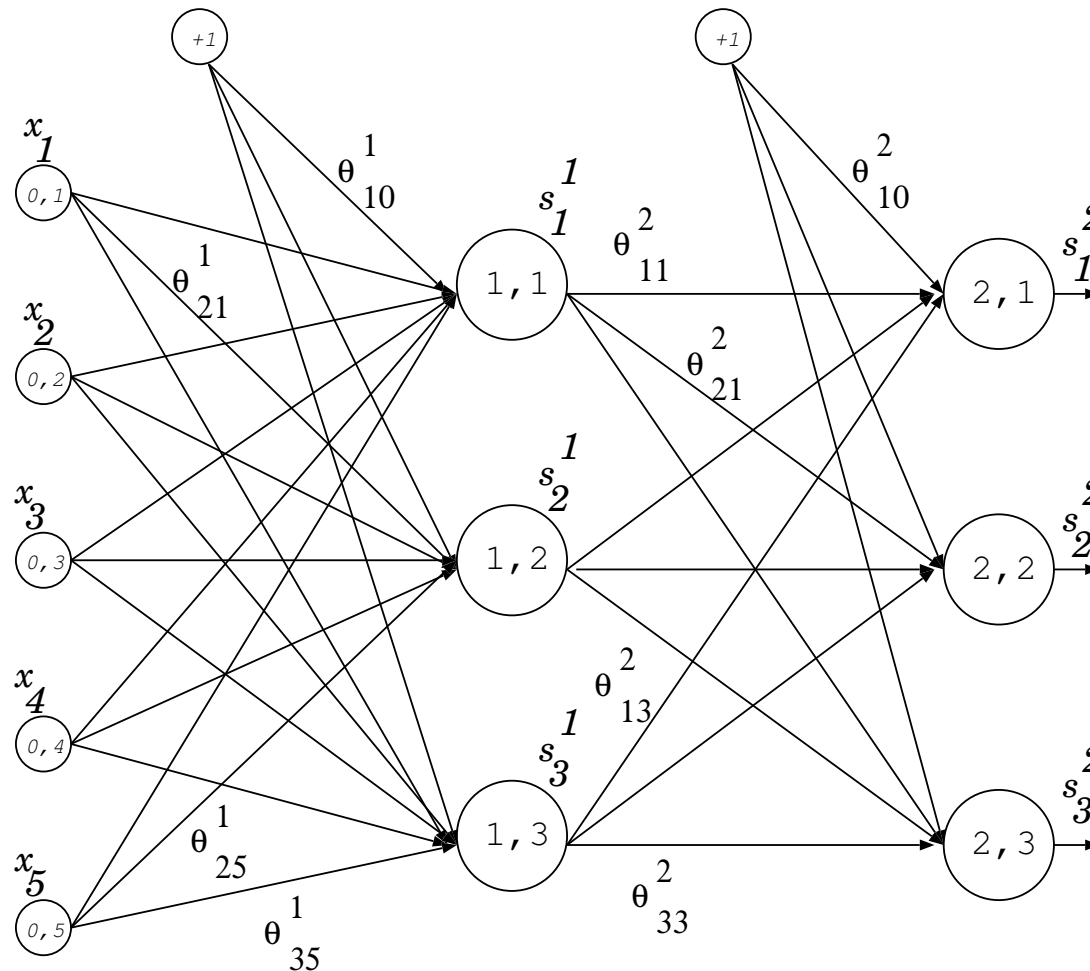
SOLUCIÓN: descenso por gradiente (“algoritmo BACKPROP”)

$$\Delta \theta_{ij}^l = -\rho \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^l} = \frac{1}{N} \sum_{n=1}^N -\rho \frac{\partial q_n(\Theta)}{\partial \theta_{ij}^l} = \frac{1}{N} \sum_{n=1}^N \Delta_n \theta_{ij}^l$$

calcular  $\Delta_n \theta_{ij}^l \stackrel{\text{def}}{=} -\rho \frac{\partial q_n(\Theta)}{\partial \theta_{ij}^l}$ ,  $1 \leq i \leq M_l$ ,  $0 \leq j \leq M_{l-1}$ ,  $1 \leq l \leq L$ ,  $1 \leq n \leq N$

Para simplificar, pero sin pérdida de generalidad, en lo que sigue se asume  $L=2$ .

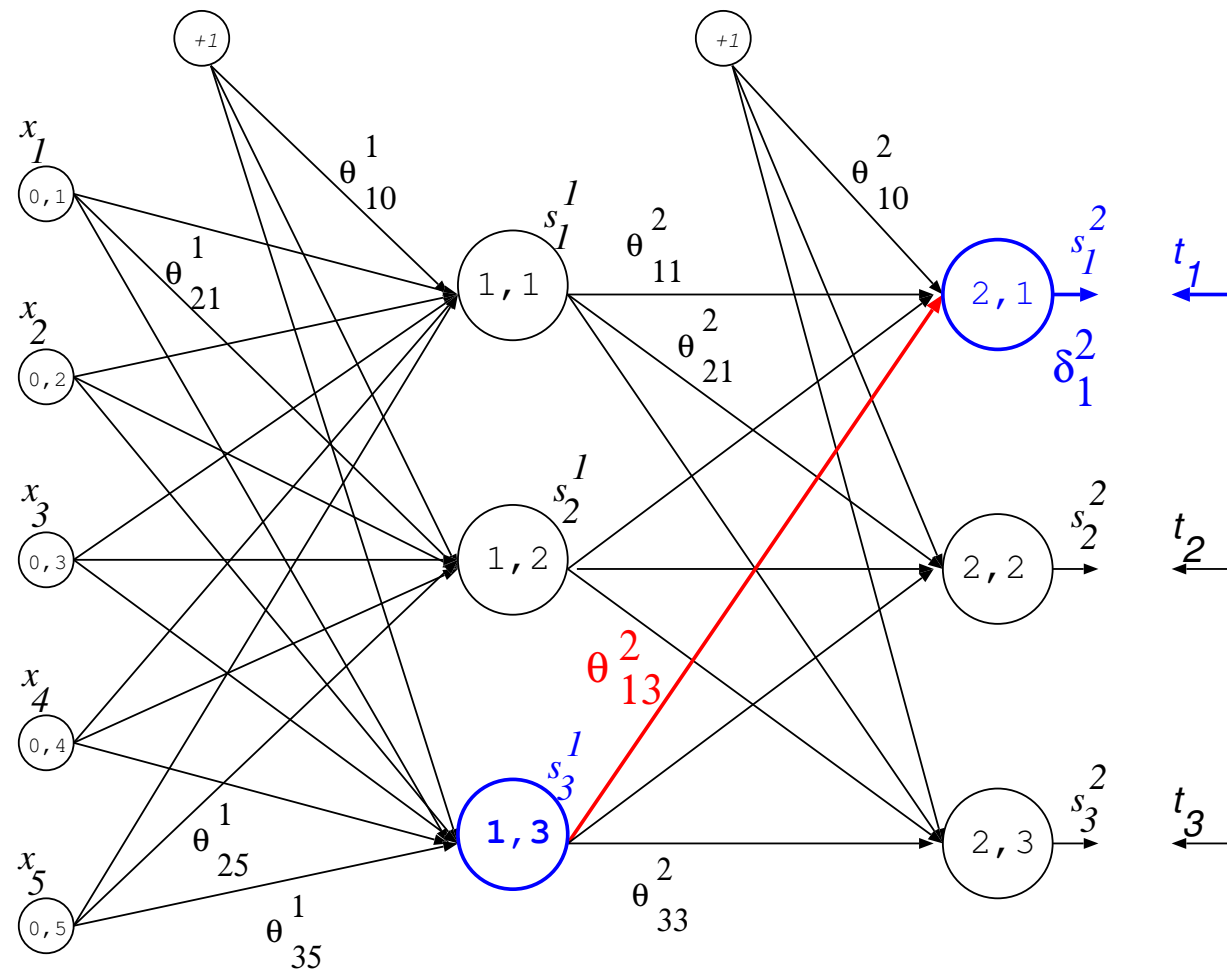
# Retropropagación del error (BackProp): cálculo hacia adelante



$$s^1_i(\mathbf{x}) = g(\phi^1_i) = g\left(\sum_{j=0}^{M_0} \theta^1_{ij} x_j\right) \quad 1 \leq i \leq M_1; \quad s^2_j = g(\phi^2_j) = g\left(\sum_{k=0}^{M_1} \theta^2_{jk} s^1_k(\mathbf{x})\right), \quad 1 \leq j \leq M_2$$

(para una muestra de entrenamiento genérica  $(\mathbf{x}, \mathbf{t})$ , y  $M_0 = 5, M_1 = 3, M_2 = 3$ )

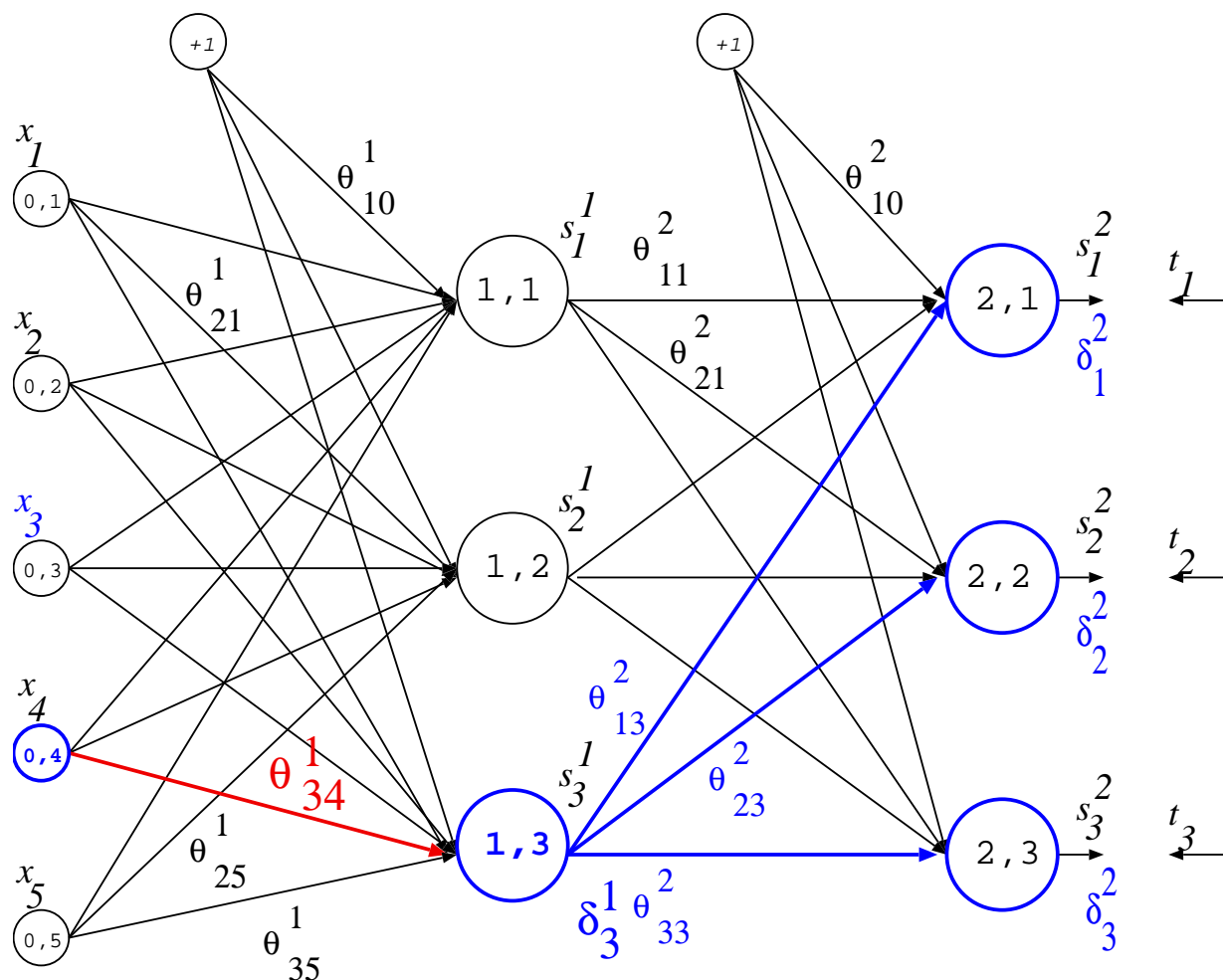
# BackProp: ilustración de actualización pesos de la capa de salida



$$\Delta \theta_{13}^2 = \rho \delta_1^2 s_3^1 = \rho (t_1 - s_1^2) g'(\phi_1^2) s_3^1$$

(para la muestra  $(x, t)$ )

# BackProp: ilustración de actualización pesos de la capa oculta



$$\Delta \theta_{34}^1 = \rho \delta_3^1 x_4 = \rho \left( g'(\phi_3^1) \sum_{r=1}^{M_2} \delta_r^2 \theta_{r3}^2 \right) x_4$$

(para la muestra  $(x, t)$ )



# Regla de la cadena para el cálculo de derivadas

- Función simple de otra función

$$f, g : \mathbb{R} \rightarrow \mathbb{R}:$$

$$\frac{d f(g(x))}{d x} = \frac{d f}{d g} \frac{d g}{d x}$$

- Función de otras dos funciones de  $n$  (o más) variables

$$g_1, g_2 : \mathbb{R}^{n \geq 2} \rightarrow \mathbb{R}, \quad f : \mathbb{R}^{M \geq 2} \rightarrow \mathbb{R}$$

$$\frac{\partial f(g_1(x, \dots), g_2(x, \dots))}{\partial x} = \frac{\partial f}{\partial g_1} \frac{\partial g_1}{\partial x} + \frac{\partial f}{\partial g_2} \frac{\partial g_2}{\partial x}$$

- Función de  $N$  (o más) funciones de  $n$  (o más) variables

$$g_1, \dots, g_N : \mathbb{R}^{n \geq N} \rightarrow \mathbb{R}, \quad f : \mathbb{R}^{M \geq N} \rightarrow \mathbb{R}:$$

$$\frac{\partial f(g_1(x, \dots), \dots, g_N(x, \dots))}{\partial x} = \sum_{i=1}^N \frac{\partial f}{\partial g_i} \frac{\partial g_i}{\partial x}$$

## Derivación del algoritmo BackProp (I)

Actualización de los pesos de la capa de salida  $\theta_{ij}^2$ , para una muestra genérica  $(\mathbf{x}, \mathbf{t}) \equiv (\mathbf{x}_n, \mathbf{t}_n)$ :

$$q(\Theta) \equiv q_n(\Theta) = \frac{1}{2} \sum_{l=1}^{M_2} (t_l - s_l^2)^2; \quad s_l^2 = g(\phi_l^2); \quad \phi_l^2 = \sum_{m=0}^{M_1} \theta_{lm}^2 s_m^1$$

$$\begin{aligned} \frac{\partial q}{\partial \theta_{ij}^2} &= \frac{\partial q}{\partial s_i^2} \frac{\partial s_i^2}{\partial \theta_{ij}^2} = \frac{\partial q}{\partial s_i^2} \frac{d s_i^2}{d \phi_i^2} \frac{\partial \phi_i^2}{\partial \theta_{ij}^2} \\ &\quad \downarrow \quad \downarrow \quad \downarrow \\ &= -((t_i - s_i^2) g'(\phi_i^2)) s_j^1 \stackrel{\text{def}}{=} -\delta_i^2 s_j^1 \end{aligned}$$

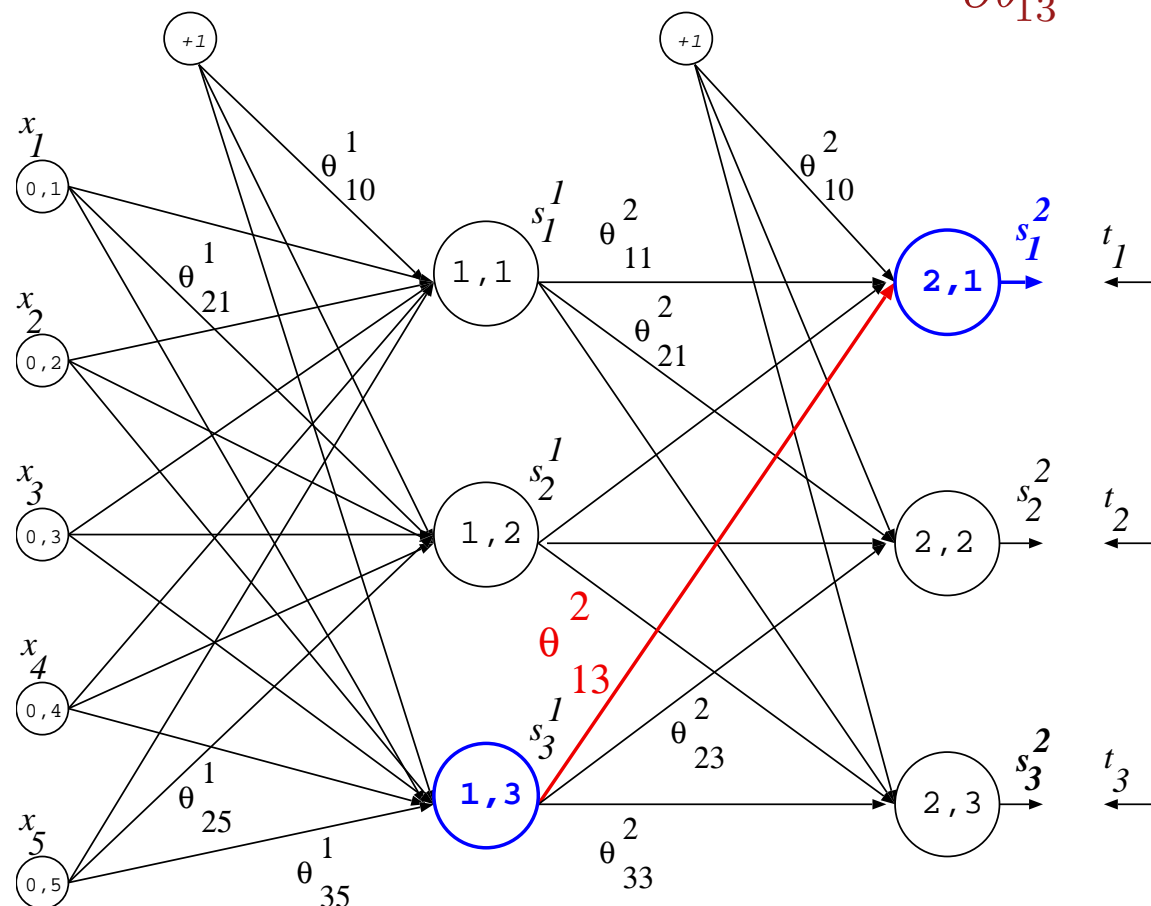
$$\frac{\partial q}{\partial \theta_{ij}^2} = -\delta_i^2 s_j^1, \quad \delta_i^2 \stackrel{\text{def}}{=} (t_i - s_i^2) g'(\phi_i^2)$$

$$\Delta_n \theta_{ij}^2 = -\rho \frac{\partial q_n}{\partial \theta_{ij}^2} = \rho \delta_i^2 s_j^1 \quad 1 \leq i \leq M_2, \quad 0 \leq j \leq M_1$$

# Ejemplo derivadas parciales respecto a pesos de la capa de salida

Para una muestra genérica  $(x, t)$ ,  $q(\Theta) = \frac{1}{2} \sum_{l=1}^3 (t_l - s_l^2(\Theta))^2$

$q(\Theta)$  depende de  $\theta_{13}^2$  solo a través de  $s_1^2(\Theta)$ :  $\frac{\partial q}{\partial \theta_{13}^2} = \frac{\partial q}{\partial s_1^2} \frac{\partial s_1^2}{\partial \theta_{13}^2}$



$$\Delta \theta_{13}^2 = \rho \delta_1^2 s_3^1 = \rho (t_1 - s_1^2) g'(\phi_1^2) s_3^1$$

## Derivación del algoritmo BackProp (II)

Actualización de los pesos de la capa oculta  $\theta_{ij}^1$ , para una muestra genérica  $(\mathbf{x}, \mathbf{t}) \equiv (\mathbf{x}_n, \mathbf{t}_n)$ :

$$q(\Theta) = \frac{1}{2} \sum_{l=1}^{M_2} (t_l - s_l^2)^2; \quad s_l^2 = g(\phi_l^2); \quad \phi_l^2 = \sum_{m=0}^{M_1} \theta_{lm}^2 s_m^1; \quad s_m^1 = g(\phi_m^1); \quad \phi_m^1 = \sum_{k=0}^{M_0} \theta_{mk}^1 x_k$$

$$\begin{aligned} \frac{\partial q}{\partial \theta_{ij}^1} &= \sum_{r=1}^{M_2} \frac{\partial q}{\partial s_r^2} \frac{\partial s_r^2}{\partial \theta_{ij}^1} = \sum_{r=1}^{M_2} \frac{\partial q}{\partial s_r^2} \frac{d s_r^2}{d \phi_r^2} \frac{\partial \phi_r^2}{\partial s_i^1} \frac{d s_i^1}{d \phi_i^1} \frac{\partial \phi_i^1}{\partial \theta_{ij}^1} \\ &= \sum_{r=1}^{M_2} \underbrace{-\delta_r^2}_{\downarrow} \underbrace{\theta_{ri}^2}_{\downarrow} \underbrace{g'(\phi_i^1)}_{\downarrow} x_j = - \left( g'(\phi_i^1) \sum_{r=1}^{M_2} \delta_r^2 \theta_{ri}^2 \right) x_j \stackrel{\text{def}}{=} -\delta_i^1 x_j \end{aligned}$$

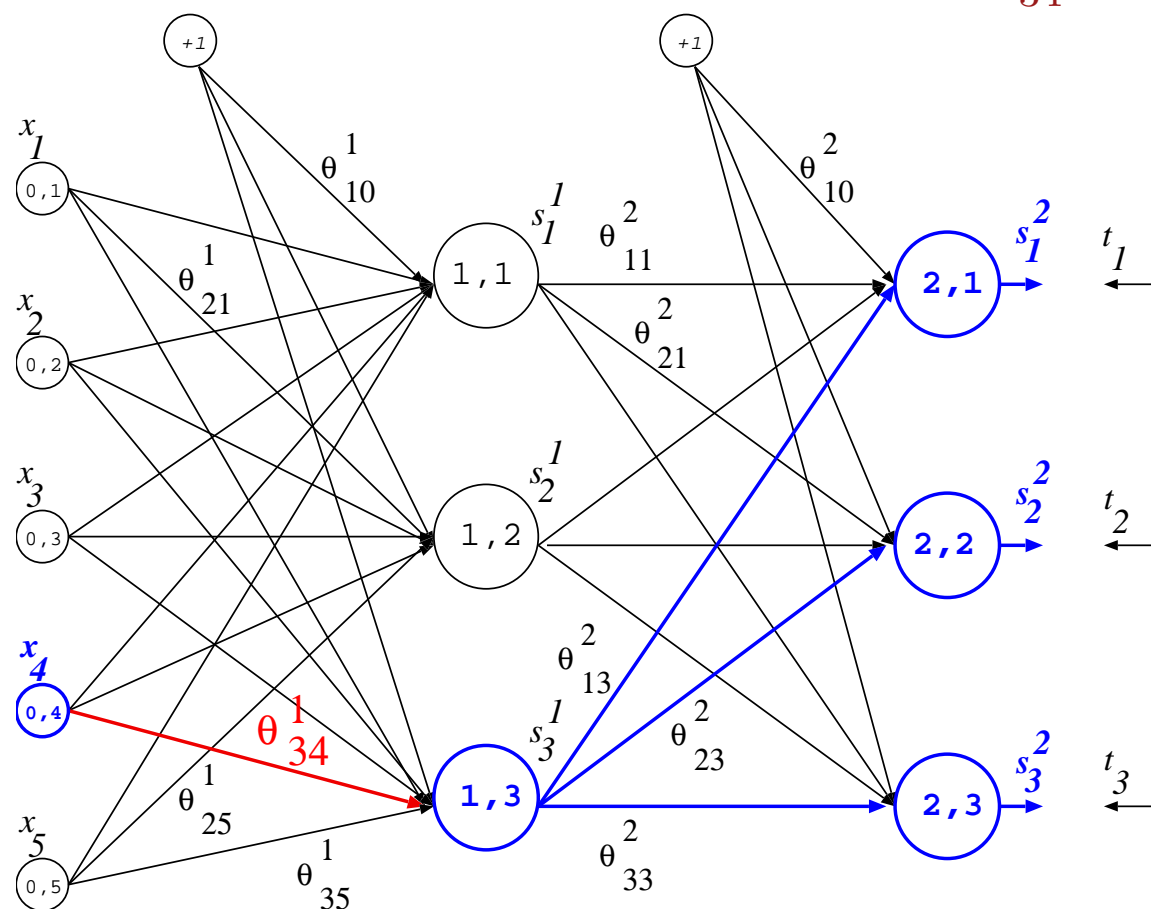
$$\frac{\partial q}{\partial \theta_{ij}^1} = -\delta_i^1 x_j, \quad \delta_i^1 \stackrel{\text{def}}{=} g'(\phi_i^1) \sum_{r=1}^{M_2} \delta_r^2 \theta_{ri}^2$$

$$\Delta_n \theta_{ij}^1 = -\rho \frac{\partial q_n}{\partial \theta_{ij}^1} = \rho \delta_i^1 x_j \quad 1 \leq i \leq M_1, \quad 0 \leq j \leq M_0$$

# Ejemplo derivadas parciales respecto a pesos de la capa oculta

Para una muestra genérica  $(x, t)$ ,  $q(\Theta) = \frac{1}{2} \sum_{l=1}^3 (t_l - s_l^2(\Theta))^2$

$q(\Theta)$  depende de  $\theta_{34}^1$  a través de  $s_r^2(\Theta)$ ,  $1 \leq r \leq 3$ :  $\frac{\partial q}{\partial \theta_{34}^1} = \sum_{r=1}^3 \frac{\partial q}{\partial s_r^2} \frac{\partial s_r^2}{\partial \theta_{34}^1}$



$$\Delta \theta_{34}^1 = \rho \delta_3^1 x_4 = \rho (g'(\phi_3^1) \sum_{r=1}^{M_2} \delta_r^2 \theta_{r3}^2) x_4$$

## Ecuaciones BackProp

- Actualización de los pesos de la capa de salida: ( $1 \leq i \leq M_2$ ,  $0 \leq j \leq M_1$ ):

$$\Delta \theta_{ij}^2 = -\rho \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^2} = \frac{\rho}{N} \sum_{n=1}^N \delta_i^2(\mathbf{x}_n) s_j^1(\mathbf{x}_n)$$

$$\delta_i^2(\mathbf{x}_n) = (t_{ni} - s_i^2(\mathbf{x}_n)) g'(\phi_i^2(\mathbf{x}_n)) \text{ con } \phi_i^2(\mathbf{x}_n) = \sum_{j=0}^{M_1} \theta_{ij}^2 s_j^1(\mathbf{x}_n)$$

- Actualización de los pesos de la capa oculta ( $1 \leq i \leq M_1$ ,  $0 \leq j \leq M_0$ ):

$$\Delta \theta_{ij}^1 = -\rho \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^1} = \frac{\rho}{N} \sum_{n=1}^N \delta_i^1(\mathbf{x}_n) x_{nj}$$

$$\delta_i^1(\mathbf{x}_n) = \left( \sum_{r=1}^{M_2} \delta_r^2(\mathbf{x}_n) \theta_{ri}^2 \right) g'(\phi_i^1(\mathbf{x}_n)) \text{ con } \phi_i^1(\mathbf{x}_n) = \sum_{j=0}^{M_0} \theta_{ij}^1 x_{nj}$$

# Ecuaciones BackProp para perceptrones de tres capas

- Actualización de los pesos de la capa de salida ( $1 \leq i \leq M_3$ ,  $0 \leq j \leq M_2$ )

$$\Delta\theta_{ij}^3 = -\rho \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^3} = \frac{\rho}{N} \sum_{n=1}^N \delta_i^3(\mathbf{x}_n) s_j^2(\mathbf{x}_n) \quad \delta_i^3(\mathbf{x}_n) = \left( t_{ni} - s_i^3(\mathbf{x}_n) \right) g'(\phi_i^3(\mathbf{x}_n))$$

- Actualización de los pesos de la segunda capa oculta ( $1 \leq i \leq M_2$ ,  $0 \leq j \leq M_1$ )

$$\Delta\theta_{ij}^2 = -\rho \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^2} = \frac{\rho}{N} \sum_{n=1}^N \delta_i^2(\mathbf{x}_n) s_j^1(\mathbf{x}_n) \quad \delta_i^2(\mathbf{x}_n) = \left( \sum_{r=1}^{M_3} \delta_r^3(\mathbf{x}_n) \theta_{ri}^3 \right) g'(\phi_i^2(\mathbf{x}_n))$$

- Actualización de los pesos de la primera capa oculta ( $1 \leq i \leq M_1$ ,  $0 \leq j \leq M_0$ )

$$\Delta\theta_{ij}^1 = -\rho \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^1} = \frac{\rho}{N} \sum_{n=1}^N \delta_i^1(\mathbf{x}_n) x_{nj} \quad \delta_i^1(\mathbf{x}_n) = \left( \sum_{r=1}^{M_2} \delta_r^2(\mathbf{x}_n) \theta_{ri}^2 \right) g'(\phi_i^1(\mathbf{x}_n))$$

# Algoritmo BACKPROP

**Entrada:** Topología, pesos iniciales  $\theta_{ij}^l$ ,  $1 \leq l \leq L$ ,  $1 \leq i \leq M_l$ ,  $0 \leq j \leq M_{l-1}$ , factor de aprendizaje  $\rho$ , condiciones de convergencia,  $N$  datos de entrenamiento  $S$

**Salidas:** Pesos de las conexiones que minimizan el error cuadrático medio de  $S$

Mientras no se cumplan las condiciones de convergencia

Para  $1 \leq l \leq L$ ,  $1 \leq i \leq M_l$ ,  $0 \leq j \leq M_{l-1}$ , inicializar  $\Delta\theta_{ij}^l = 0$

Para cada muestra de entrenamiento  $(x, t) \in S$

Desde la capa de entrada a la de salida ( $l = 0, \dots, L$ ):

Para  $1 \leq i \leq M_l$  si  $l = 0$  entonces  $s_i^0 = x_i$  sino calcular  $\phi_i^l$  y  $s_i^l = g(\phi_i^l)$

Desde la capa de salida a la de entrada ( $l = L, \dots, 1$ ),

Para cada nodo ( $1 \leq i \leq M_l$ )

Calcular  $\delta_i^l = \begin{cases} g'(\phi_i^l) (t_{ni} - s_i^L) & \text{si } l == L \\ g'(\phi_i^l) (\sum_r \delta_r^{l+1} \theta_{ri}^{l+1}) & \text{en otro caso} \end{cases}$

Para cada peso  $\theta_{ij}^l$  ( $0 \leq j \leq M_{l-1}$ ) calcular:  $\Delta\theta_{ij}^l = \Delta\theta_{ij}^l + \rho \delta_i^l s_j^{l-1}$

Para  $1 \leq l \leq L$ ,  $1 \leq i \leq M_l$ ,  $0 \leq j \leq M_{l-1}$ , actualizar pesos:  $\theta_{ij}^l = \theta_{ij}^l + \frac{1}{N} \Delta\theta_{ij}^l$

Coste computacional por cada iteración *mientras*:  $O(N D)$ ,  $N = |S|$ ,  $D$  = número de pesos

DEMO: <http://playground.tensorflow.org/>



## Algoritmo BACKPROP (“incremental”)

**Entrada:** Topología, pesos iniciales  $\theta_{ij}^l$ ,  $1 \leq l \leq L$ ,  $1 \leq i \leq M_l$ ,  $0 \leq j \leq M_{l-1}$ ,  
factor de aprendizaje  $\rho$ , condiciones de convergencia,  $N$  datos de entrenamiento  $S$

**Salidas:** Pesos de las conexiones que minimizan el error cuadrático medio de  $S$

Mientras no se cumplan las condiciones de convergencia

Para cada muestra de entrenamiento  $(x, t) \in S$  (en orden aleatorio)

Desde la capa de entrada a la de salida ( $l = 0, \dots, L$ ):

Para  $1 \leq i \leq M_l$  si  $l = 0$  entonces  $s_i^0 = x_i$  sino calcular  $\phi_i^l$  y  $s_i^l = g(\phi_i^l)$

Desde la capa de salida a la de entrada ( $l = L, \dots, 1$ ),

Para cada nodo ( $1 \leq i \leq M_l$ )

Calcular  $\delta_i^l = \begin{cases} g'(\phi_i^l) (t_{ni} - s_i^L) & \text{si } l == L \\ g'(\phi_i^l) (\sum_r \delta_r^{l+1} \theta_{ri}^{l+1}) & \text{en otro caso} \end{cases}$

Para cada peso  $\theta_{ij}^l$  ( $0 \leq j \leq M_{l-1}$ ) calcular:  $\Delta\theta_{ij}^l = \rho \delta_i^l s_j^{l-1}$

Para  $1 \leq l \leq L$ ,  $1 \leq i \leq M_l$ ,  $0 \leq j \leq M_{l-1}$ , actualizar pesos:  $\theta_{ij}^l = \theta_{ij}^l + \frac{1}{N} \Delta\theta_{ij}^l$

Coste computacional por cada iteración *mientras*:  $O(N D)$ ,  $N = |S|$ ,  $D =$  número de pesos

## Algoritmo BACKPROP (“on-line”)

**Entrada:** Topología, pesos iniciales  $\theta_{ij}^l$ ,  $1 \leq l \leq L$ ,  $1 \leq i \leq M_l$ ,  $0 \leq j \leq M_{l-1}$ ,  
factor de aprendizaje  $\rho$ , dato de entrenamiento  $(x, t)$

**Salidas:** Pesos de las conexiones actualizados mediante  $(x, t)$

Desde la capa de entrada a la de salida ( $l = 0, \dots, L$ ):

Para  $1 \leq i \leq M_l$  si  $l = 0$  entonces  $s_i^0 = x_i$  sino calcular  $\phi_i^l$  y  $s_i^l = g(\phi_i^l)$

Desde la capa de salida a la de entrada ( $l = L, \dots, 1$ ),

Para cada nodo ( $1 \leq i \leq M_l$ )

Calcular  $\delta_i^l = \begin{cases} g'(\phi_i^l) (t_{ni} - s_i^L) & \text{si } l == L \\ g'(\phi_i^l) (\sum_r \delta_r^{l+1} \theta_{ri}^{l+1}) & \text{en otro caso} \end{cases}$

Para cada peso  $\theta_{ij}^l$  ( $0 \leq j \leq M_{l-1}$ ) calcular:  $\Delta\theta_{ij}^l = \rho \delta_i^l s_j^{l-1}$

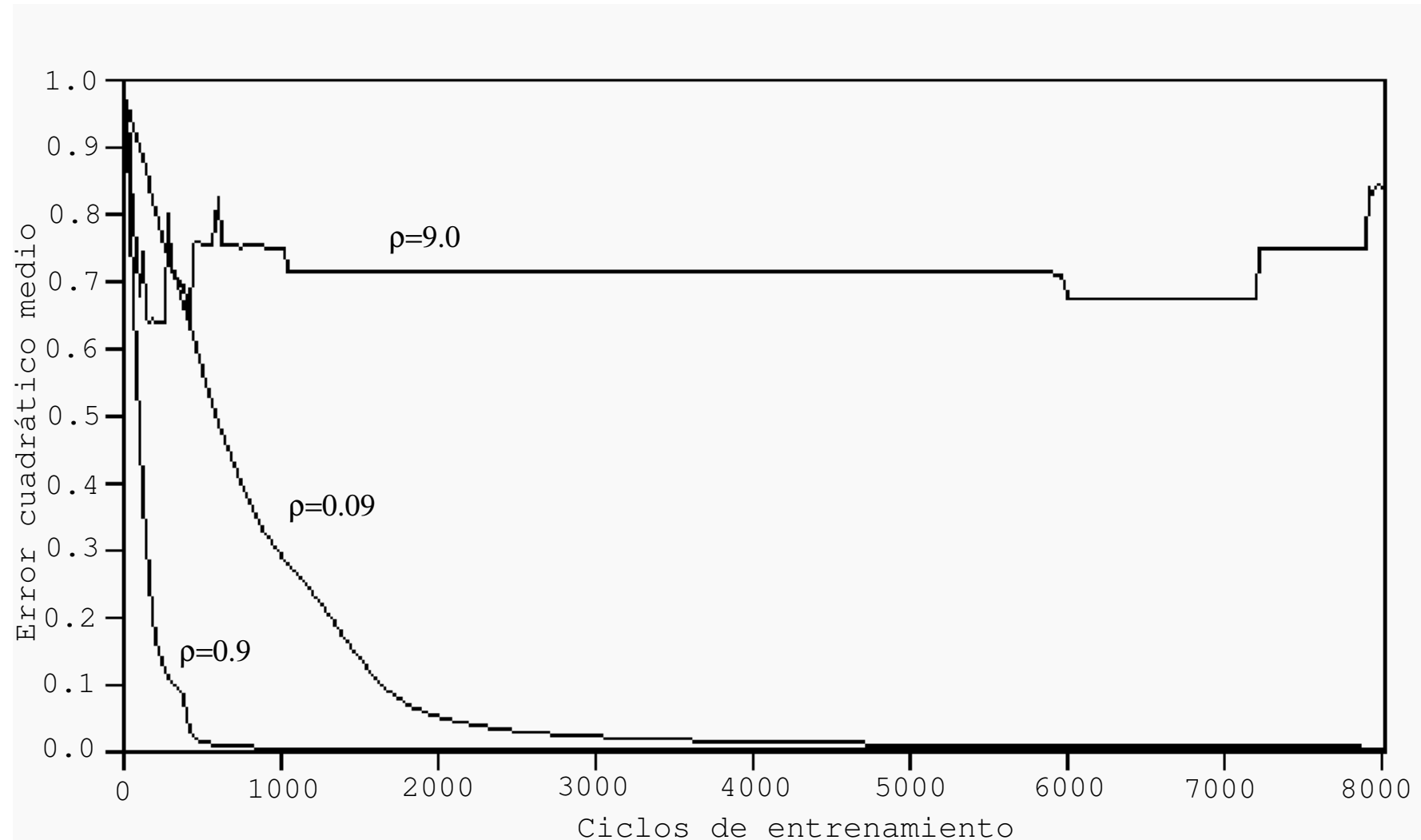
Para  $1 \leq l \leq L$ ,  $1 \leq i \leq M_l$ ,  $0 \leq j \leq M_{l-1}$ , actualizar pesos:  $\theta_{ij}^l = \theta_{ij}^l + \Delta\theta_{ij}^l$

Coste computacional por cada muestra procesada:  $O(D)$ ,  $D = \text{número de pesos}$

# Index

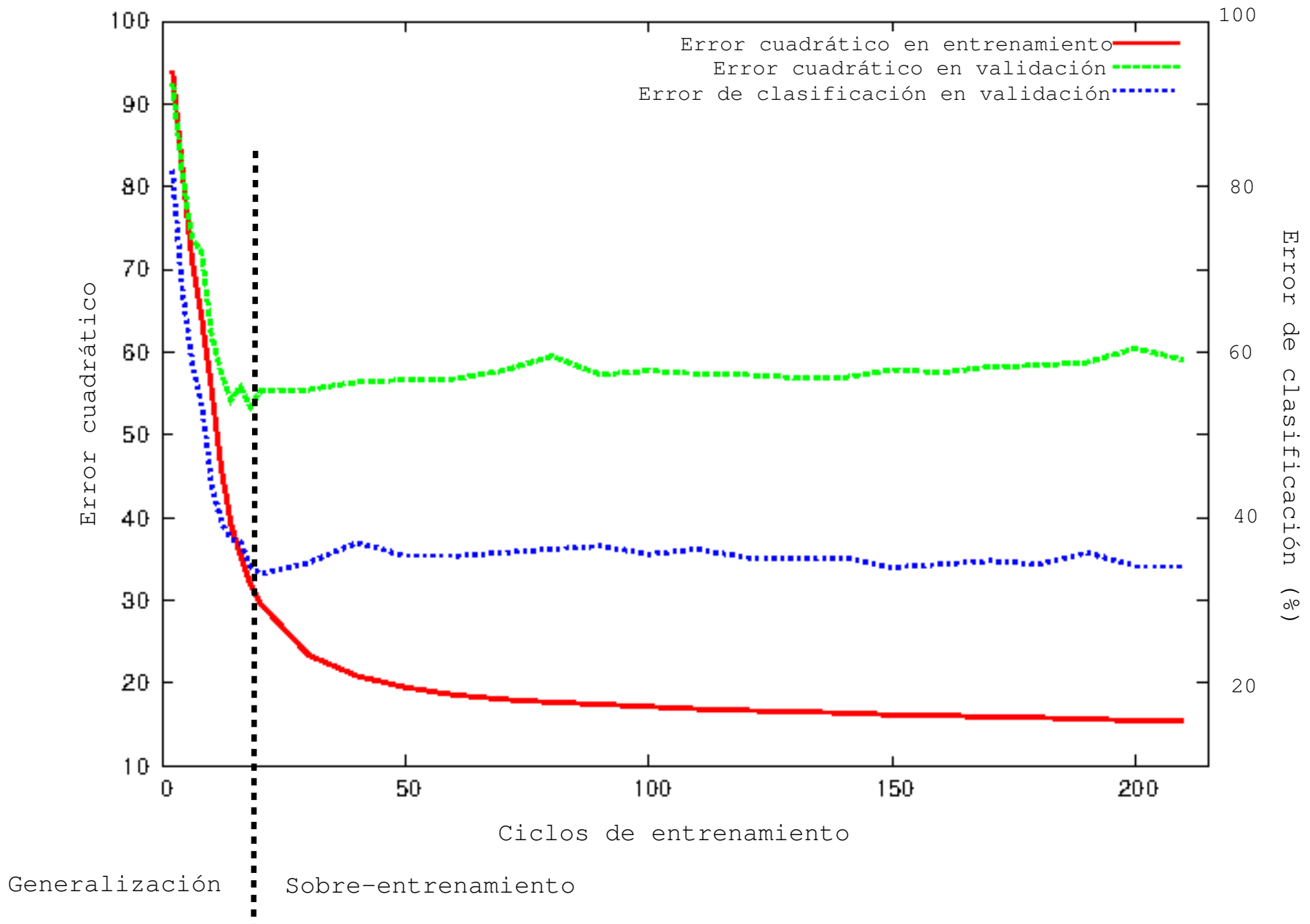
- 1 Redes neuronales multicapa ▷ 2
- 2 Algoritmo de retropropagación del error (BackProp) ▷ 18
- 3 *Aspectos de uso y propiedades del BackProp* ▷ 33
- 4 Variantes de BackProp ▷ 41
- 5 Redes neuronales radiales ▷ 47
- 6 Aplicaciones ▷ 52
- 7 Notación ▷ 54

# Selección del factor de aprendizaje



<http://playground.tensorflow.org/>

# Condiciones de convergencia



## Algunos problemas y soluciones con el BackProp

- El problema de la anulación o explosión del gradiente en el caso de muchas capas:
  - Normalización de la entrada y normalización por capa.
  - Inicialización de los pesos a valores pequeños.
  - Congelar algunos pesos de forma aleatoria (Dropout).
  - Recorte del gradiente (Adaptive gradient clipping).
  - Conexiones residuales.
- Evitar “malos” mínimos locales:
  - Barajar los datos de entrenamiento.
  - Aprender primero las muestras más “fáciles” (Curriculum learning).
  - Regularización.
  - Añadir ruido durante el entrenamiento.

## Codificación

- Los valores  $\pm 1$  (o  $0, 1$ ) solo se alcanzan asintóticamente cuando se utilizan la mayoría de las funciones de activación como la sigmoid:

Valores deseados de salida:  $t_i = \begin{cases} -0.9 \\ +0.9 \end{cases}$

- Parálisis de la red: para valores grandes de  $z$  la derivada  $g'(z)$  es muy pequeña y por tanto, los incrementos de los pesos son muy pequeños. Una forma de disminuir este efecto es **normalizar el rango de entrada**.

$$S = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subset \mathbb{R}^d \Rightarrow \begin{cases} \mu_j = \frac{1}{N} \sum_{i=1}^N x_{ij} & 1 \leq j \leq d \\ \sigma_j^2 = \frac{1}{N-1} \sum_{i=1}^N (x_{ij} - \mu_j)^2 & 1 \leq j \leq d \end{cases}$$

$$\forall \mathbf{x} \in \mathbb{R}^d, \hat{\mathbf{x}} : \hat{x}_j = \frac{x_j - \mu_j}{\sigma_j} \Rightarrow \begin{cases} \hat{\mu}_j = 0 \\ \hat{\sigma}_j = 1 \end{cases} \text{ for } 1 \leq j \leq d$$

## Propiedades del BackProp

- **Convergencia:** teorema general del descenso por gradiente (Tema 3)
- **Elección del factor de aprendizaje:** Tema 3 (Adadelata, Adam, ...)
- **Coste computacional:**  $O(N D)$  en cada iteración
- **Perceptrones de una o dos capas ocultas** [Villiams and Basnard, 92]
  - En el mejor de los casos, los resultados de clasificación no presentan diferencias estadísticas significativas.
  - Generalmente, un perceptrón de una capa oculta suele producir mejores resultados de clasificación que los de dos capas ocultas.
  - Los perceptrones de dos capas ocultas suelen converger más rápidamente que los perceptrones de una capa oculta.
- *En condiciones límites, las salidas de un perceptrón entrenado para minimizar el error cuadrático medio de las muestras de entrenamiento de un problema de clasificación aproximan la distribución a-posteriori subyacente en las muestras de entrenamiento.*



# Estimación de la topología

- MÉTODOS DE PODA:
  - Poda de pesos basadas en:
    - \* RELEVANCIA
    - \* CASTIGO
  - Poda de FDLA.
- MÉTODOS INCREMENTALES:
  - Búsqueda incremental (Moddy & Utans, 1995)
  - “Cascade-correlation” (Fahlman & Lebiens, 1990)
  - Adaptación estructural (Lee et al. 1990)
- POR TRANSFORMACIÓN
  - Árboles de decisión

## Aspectos computacionales del aprendizaje con el PM

- La dimension de Vapnik-Chervonenkis (transparencias 3-9 del tema 2) de un perceptrón multicapa es el número de conexiones  $D$ .
- Sobre el tamaño del conjunto de entrenamiento (Ripley, 1993 y tema 2.1): Si el perceptrón multicapa tiene  $M$  nodos,  $D$  pesos y un número arbitrario de capas, el número de muestras de entrenamiento debe ser proporcional a  $\frac{D}{\epsilon} \log \frac{M}{\epsilon}$  donde  $\epsilon$  es el error tolerado. Si el perceptrón multicapa tiene solo una capa, el número de muestras de entrenamiento debe ser proporcional a  $\frac{D}{\epsilon}$ .

# Index

- 1 Redes neuronales multicapa ▷ 2
- 2 Algoritmo de retropropagación del error (BackProp) ▷ 18
- 3 Aspectos de uso y propiedades del BackProp ▷ 33
- 4 *Variantes de BackProp* ▷ 41
- 5 Redes neuronales radiales ▷ 47
- 6 Aplicaciones ▷ 52
- 7 Notación ▷ 54

## BackProp específico para clasificación

- En problemas de clasificación se suele utilizar la función de activación *softmax* en la capa de salida. En este caso, las salidas de la red pueden considerarse como aproximaciones a las probabilidades a posteriori de clase.
- Dado un conjunto de entrenamiento  $S = \{(\mathbf{x}_1, \mathbf{t}_1), \dots, (\mathbf{x}_N, \mathbf{t}_N)\}$ , con  $\mathbf{x}_n \in \mathbb{R}^{M_0}$ ,  $\mathbf{t}_n \in \{0, 1\}^{M_2}$ , ( $M_2 \equiv C$ ), esto permite establecer como criterio de optimización, alternativo al error cuadrático, la **entropía cruzada**:

$$q_S(\Theta) = -\frac{1}{N} \sum_{n=1}^N \sum_{i=1}^{M_2} t_{ni} \log s_i^2(\mathbf{x}_n; \Theta)$$

- *Problema de entrenamiento*: encontrar  $\Theta$  tal que la **entropía cruzada** sea mínima.  
*Solución*: **DESCENSO POR GRADIENTE**:

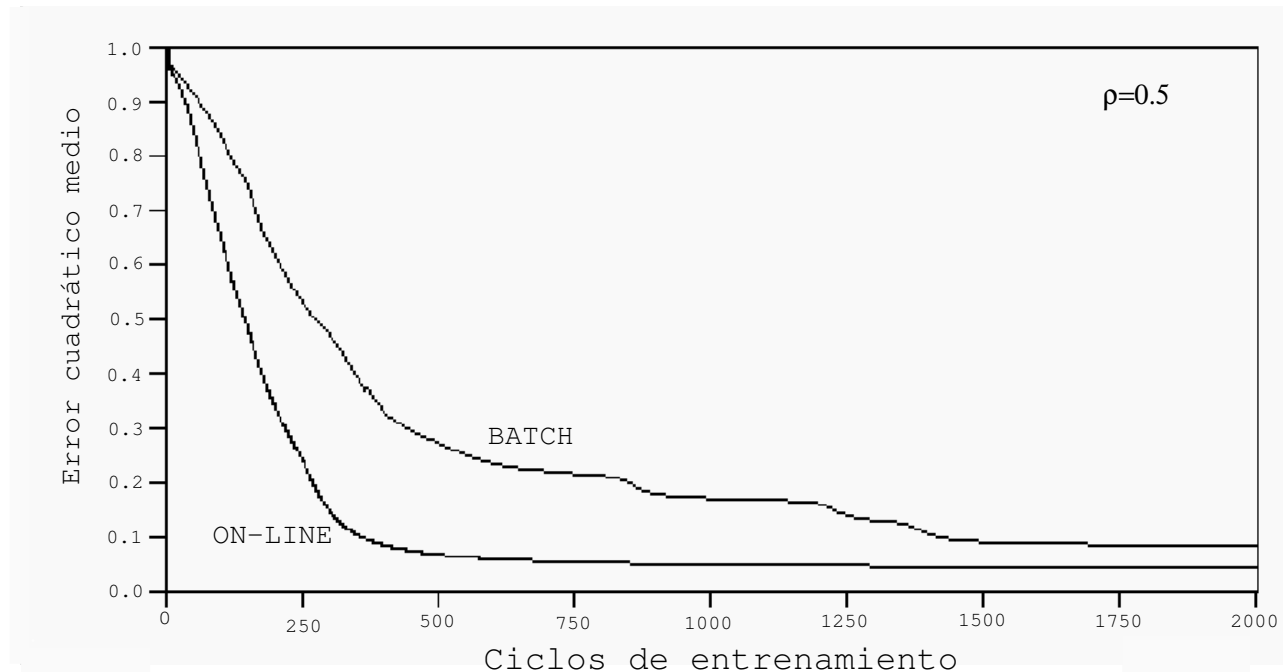
$$\Delta \theta_{ij}^l = -\rho \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^l} \quad 1 \leq l \leq 2, \quad 1 \leq i \leq M_l, \quad 0 \leq j \leq M_{l-1}$$

- El algoritmo basado en la minimización de la entropía cruzada suele producir entrenamientos más rápidos y mejores generalizaciones (Bishop 2006)

**Ejercicio:** Obtener las ecuaciones de actualización para la minimización de la entropía cruzada.

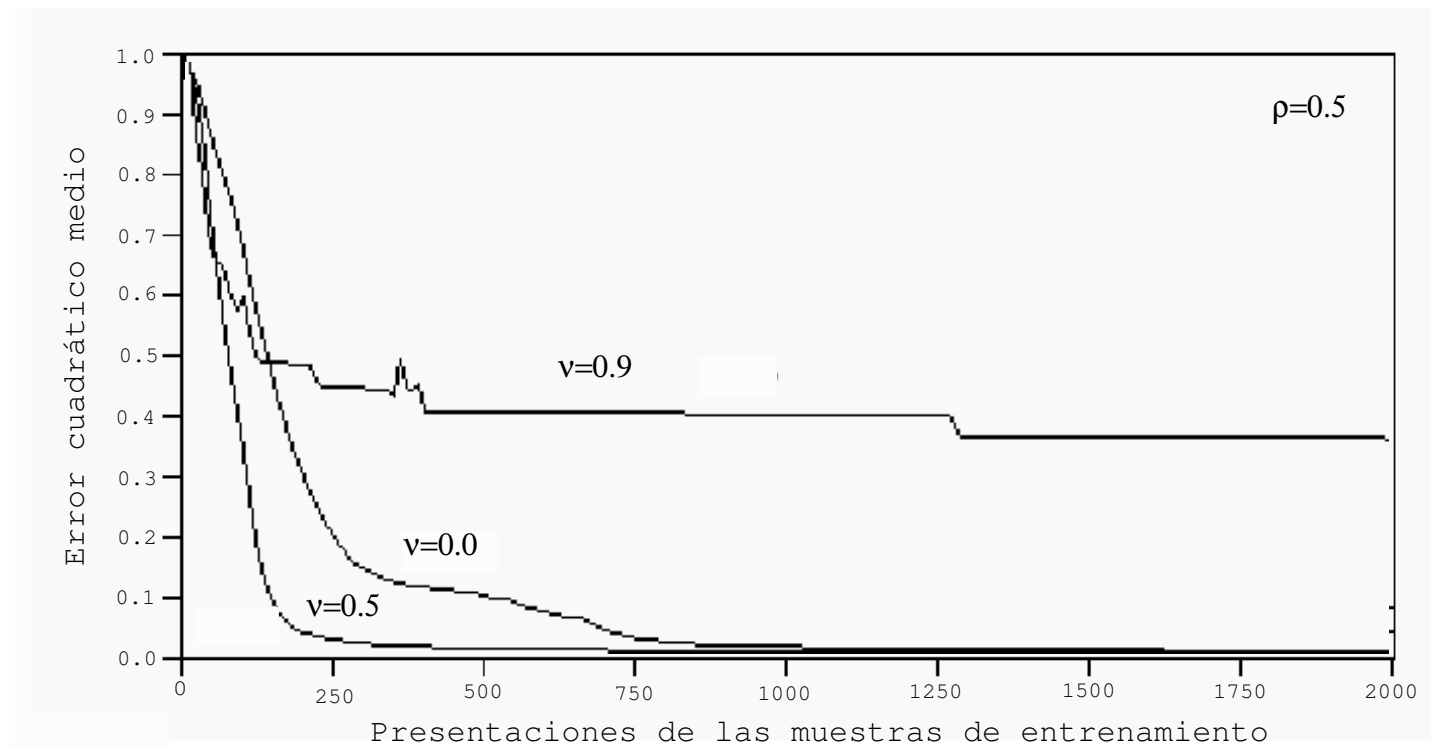
## BackProp incremental o “batch”

- BackProp “batch”: en cada iteración se procesan las  $N$  muestras de entrenamiento (“epoch”) y los pesos del PM se actualizan una sola vez.
- BackProp “*mini-batch*”: el conjunto de entrenamiento se divide en  $B$  bloques, en cada uno se procesan las muestras que están contenidas y luego se actualizan los pesos del PM. Por tanto en un epoch los pesos se actualizan  $N/B$  veces.
- BackProp “*incremental*”: en cada iteración se procesa solo una muestra de entrenamiento (aleatoria) y se actualizan los pesos del PM. Por tanto en un epoch los pesos se actualizan  $N$  veces.



## BackProp con momentum

- Problema: convergencia lenta en “plateaux”
- Posible solución: añadir una “inercia” o “momentum” con un peso  $0 \leq \nu < 1$
- BP con momentum (batch): 
$$\Delta\theta_{ij}^l(k) = \frac{\rho_k}{N} \sum_{n=1}^N \delta_i^l(\mathbf{x}_n) s_j^{l-1}(\mathbf{x}_n) + \nu \Delta\theta_{ij}^l(k-1)$$
- **TEOREMA** (Phansalkar and Sartry, 1994): *Los puntos estables del algoritmo BackProp con momentum ( $\Theta(k) = \Theta(k+1)$ ) son mínimos locales de  $q_S(\Theta)$*



## BackProp con amortiguamiento

- Problema: evitar que los pesos sean muy grandes y provoquen una parálisis de la red.
- Solución: minimizar el error cuadrático medio con **regularización**:

$$q_S(\Theta) + \frac{\lambda}{2} \sum_{l,i,j} (\theta_{ij}^l)^2$$

- Algoritmo “batch”, para  $1 \leq l \leq L$ ,  $1 \leq i \leq M_l$ ,  $1 \leq j \leq M_{l-1}$

$$\Delta \theta_{ij}^l(k) = \frac{\rho_k}{N} \sum_{n=1}^N \delta_i^l(\mathbf{x}_n) s_j^{l-1}(\mathbf{x}_n) - \rho_k \lambda \theta_{ij}^l(k-1)$$

# Algoritmo de Newton

$$\Theta(k+1) = \Theta(k) - H^{-1} \nabla q_S$$

donde  $H$  es la matrix de Hessian:

$$H_{lij,mrs} = \frac{\partial^2 q_S(\Theta)}{\partial \theta_{ij}^l \partial \theta_{rs}^m}$$

- Calcular la matrix de Hessian (a partir de las derivadas del BackProp):
  - (Ripley, 1996) pp. 151-153: capa oculta.
  - (Bishop, 1995) pp. 150-160: general, aproximación diagonal, ...
  - Una aproximación iterativa: Algoritmo de Quasi-Newton.
- Otras aplicaciones de la matrix de Hessian *métodos de poda de la red neuronal*



# Index

- 1 Redes neuronales multicapa ▷ 2
- 2 Algoritmo de retropropagación del error (BackProp) ▷ 18
- 3 Aspectos de uso y propiedades del BackProp ▷ 33
- 4 Variantes de BackProp ▷ 41
- 5 *Redes neuronales radiales* ▷ 47
- 6 Aplicaciones ▷ 52
- 7 Notación ▷ 54

## Redes neuronales radiales

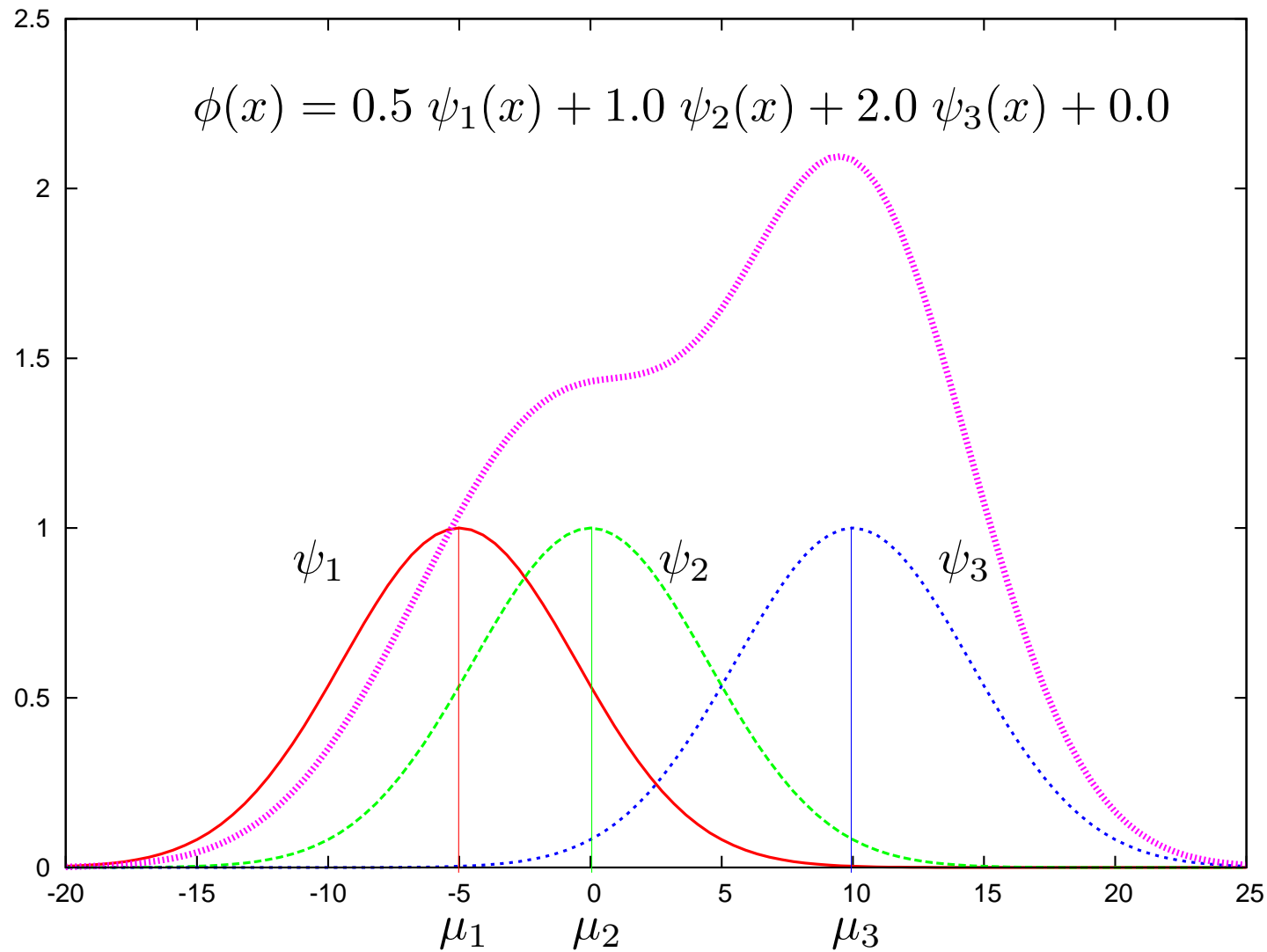
- *Una red neuronal radial* (RNR) es una FDLG de la forma:

$$\phi_m(\mathbf{x}; \Theta) = \sum_{i=1}^{d'} \theta_{mi} \psi(\|\mathbf{x} - \boldsymbol{\mu}_i\|) + \theta_{m0} \quad 1 \leq m \leq M$$

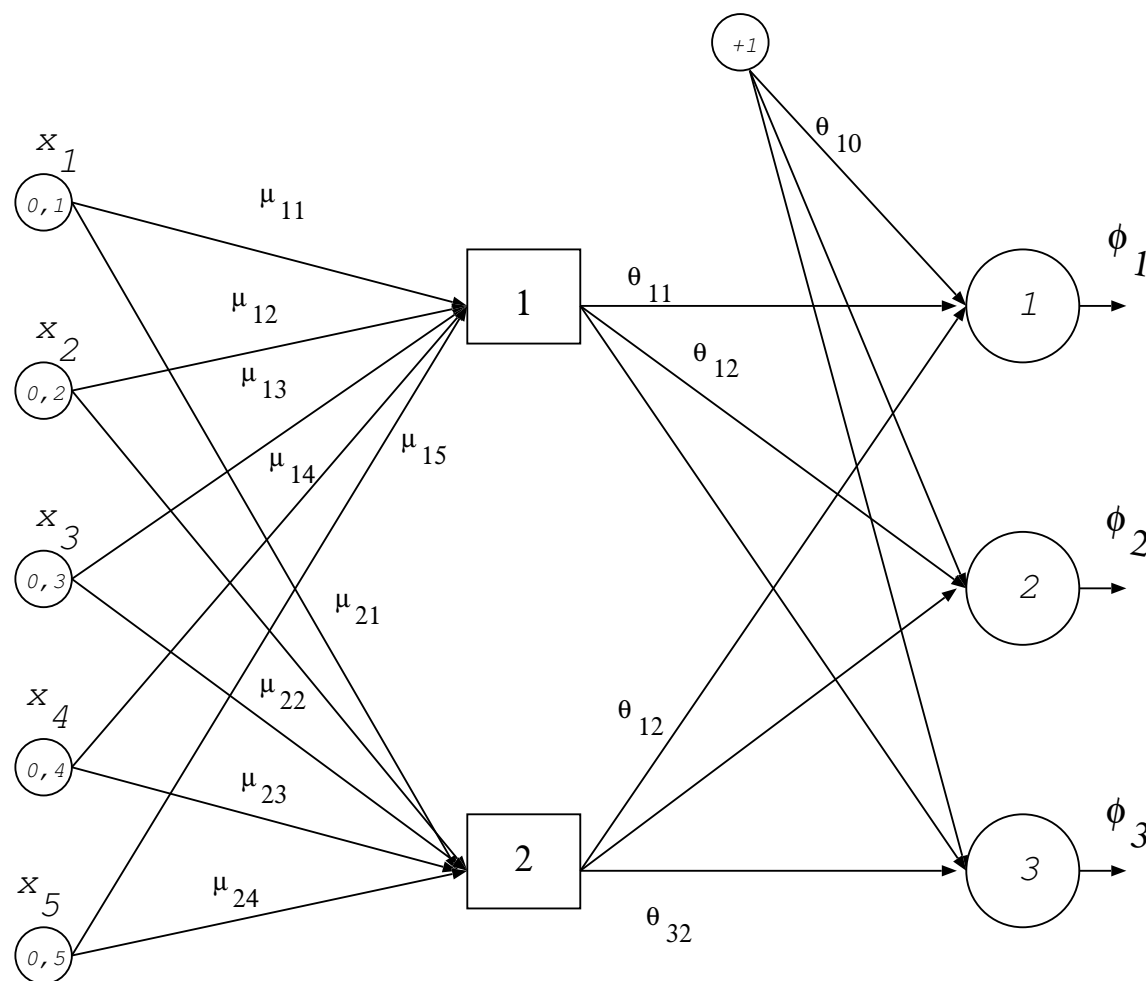
donde:  $\mathbf{x} \in \mathbb{R}^d$ ,  $\boldsymbol{\mu}_i \in \mathbb{R}^d$ , para  $1 \leq i \leq d'$   $\boldsymbol{\theta}_m \in \mathbb{R}^{d'}$ ,  $\theta_{m0} \in \mathbb{R}$  y  $\psi$  es una **función básica radial** típicamente de la forma:  $\psi(z) = \exp\left(-\frac{z^2}{2\sigma^2}\right)$  con  $\sigma \in \mathbb{R}$ .

- Para clasificación,  $M \equiv C$ , aunque las RNR son populares para regresión.
- Aprendizaje de RNR:
  - **Aprendizaje secuencial**, primero de las funciones básicas radiales (clustering) y a continuación los pesos (Adaline, ...).
  - **Aprendizaje integrado** de las funciones básicas radiales y los pesos mediante la minimización del error cuadrático medio (similar al BackProp).

## Ejemplo de funciones básicas radiales en $\mathbb{R}$



# Redes neuronales radiales



$$1 \leq i \leq d = 5$$

$$1 \leq i \leq d' = 2$$

$$1 \leq i \leq M = 3$$

$$x_i \in \mathbb{R} \quad \psi(\mathbf{x}; \boldsymbol{\mu}_i, \sigma) = \exp \left( -\frac{\sum_{j=1}^d (\mu_{ij} - x_j)^2}{2\sigma^2} \right) \quad \phi_i(\mathbf{x}; \boldsymbol{\Theta}) = \sum_{j=1}^{d'} \theta_{ij} \psi(\mathbf{x}; \boldsymbol{\mu}_j, \sigma) + \theta_{i0}$$

# Aprendizaje secuencial de las redes neuronales radiales

Dado un conjunto de entrenamiento  $S = \{(\mathbf{x}_1, \mathbf{t}_1), \dots, (\mathbf{x}_N, \mathbf{t}_N)\}$ , con  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $\mathbf{t}_i \in \mathbb{R}^M$ ,  $1 \leq i \leq N$

- **Aprendizaje secuencial** de las funciones radiales  $\psi(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\sigma}_k)$  para  $1 \leq k \leq d'$  y los pesos  $\boldsymbol{\theta}_m \in \mathbb{R}^{d'+1}$  para  $1 \leq m \leq M$ :
  1. Aprendizaje de las funciones radiales a partir de  $A' = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ ,
    - usando *clustering* (SIN)
    - usando *mixtures de gaussianas*
    - usando *mapas autoorganizativos*
  2. Aprendizaje de los pesos:  $A'' = \{(\psi(\mathbf{x}_1), \mathbf{t}_1), \dots, (\psi(\mathbf{x}_n), \mathbf{t}_n)\}$ ,
    - Perceptron (SIN y APR)
    - Pocket Perceptron
    - Widrow-Hoff (APR)

# Index

- 1 Redes neuronales multicapa ▷ 2
- 2 Algoritmo de retropropagación del error (BackProp) ▷ 18
- 3 Aspectos de uso y propiedades del BackProp ▷ 33
- 4 Variantes de BackProp ▷ 41
- 5 Redes neuronales radiales ▷ 47
- 6 *Aplicaciones* ▷ 52
- 7 Notación ▷ 54

## Algunas aplicaciones

- CLASIFICACIÓN

- Reconocimiento de caracteres impresos y manuscritos
- Exploración petrolífera.
- Aplicaciones médicas: *detección de ataques de epilepsia, ayuda al diagnóstico de la esclerosis múltiple*, etc.
- Minería de datos

- PREDICCIÓN

- Previsión de ocupación en los vuelos (Inc. BehavHeuristics.)
- Evolución de los precios de solares (Ayuntamiento de Boston)
- Predicción de consumo de bebidas refrescantes (Britvic)
- Predicción meteorológica (National Weather Service)
- Predicción de stocks (Carl & Associates, Neural Applications Corporation, etc.)
- Predicción de demanda eléctrica (Bayernwerk AG, Britvic, etc.)
- Predicción de fallos en motores eléctricos (Siemens)

- CONTROL AND AUTOMATIZATION

- Refinado del petróleo (Texaco).
- Producción de acero (Fujitsu, Neural Applications Corporation, Nippon Steel)

- + APLICACIONES <http://www.calsci.com/Applications.html>

# Index

- 1 Redes neuronales multicapa ▷ 2
- 2 Algoritmo de retropropagación del error (BackProp) ▷ 18
- 3 Aspectos de uso y propiedades del BackProp ▷ 33
- 4 Variantes de BackProp ▷ 41
- 5 Redes neuronales radiales ▷ 47
- 6 Aplicaciones ▷ 52
- 7 *Notación* ▷ 54



## Notación

- **Funciones discriminantes lineales:**  $\phi(\mathbf{x}; \Theta) = \Theta^t \mathbf{x}$  para una entrada  $\mathbf{x}$  y parámetros  $\Theta$  compuestos por vector de pesos y umbral  $(\theta, \theta_0)$
- **Funciones discriminantes lineales con activación:**  $g \circ \phi(\mathbf{x}; \theta)$  para una entrada  $\mathbf{x}$ , parámetros  $(\theta, \theta_0)$  y  $g$  una función de activación.  $g'$  es la derivada de la función de activación  $g$
- **Función de activación sigmoid:**  $g_S(z)$
- **Salida del nodo  $i$  en la capa  $k$ :**  $s_i^k$  en perceptrones multicapa y redes hacia adelante
- **Pesos de la conexión** que va del nodo  $j$  de la capa  $k - 1$  al nodo  $i$  de la capa  $k$  en un perceptrón multicapa:  $\theta_{ij}^k$ .  $\Theta$  es un vector de talla  $D$  formado por todos los pesos  $\theta_{ij}^k$ . Pesos de la conexión que va del nodo  $j$  de la capa  $k'$  al nodo  $i$  de la capa  $k$  en una red hacia adelante:  $\theta_{ij}^{k',k}$
- **Conjunto de  $N$  muestras de entrenamiento:**  $S = \{(\mathbf{x}_1, \mathbf{t}_1), \dots, (\mathbf{x}_N, \mathbf{t}_N)\}$  con  $\mathbf{x}_n \in \mathbb{R}^{M_0}$  y  $\mathbf{t}_n \in \mathbb{R}^{M_K}$ , siendo  $K$  el número de capas y  $M_k$  el número de nodos de la capa  $k$
- **Función a minimizar** en el entrenamiento de un perceptrón multicapa:  $q_S(\Theta) \in \mathbb{R}$
- **Clasificador** en  $C \equiv M_K$  clases de puntos de  $\mathbb{R}^d \equiv \mathbb{R}^{M_0}$ :  $f : \mathbb{R}^{M_0} \rightarrow \{1, \dots, M_K\}$
- **Error en el nodo  $i$  de la capa  $k$  para la muestra  $\mathbf{x}_n$ :**  $\delta_i^k(\mathbf{x}_n)$
- **Incremento del peso** que va del nodo  $j$  en la capa  $k - 1$  al nodo  $i$  en la capa  $k$ :  $\Delta \theta_{ij}^k$
- **Factor de aprendizaje, momentum y factor de regularización:**  $\rho$ ,  $\nu$  y  $\lambda$
- **Media y desviación típica:**  $\mu$  y  $\sigma$