

2016-2017

Aprendizaje Automático

6. 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)

Septiembre, 2016

Página intencionadamente en blanco

Index

- 1 *Redes neuronales multicapa* ▷ 2
- 2 Algoritmo de retropropagación del error (BackProp) ▷ 18
- 3 Variantes del BackProp ▷ 35
- 4 Redes neuronales radiales ▷ 39
- 5 Aplicaciones ▷ 43
- 6 Notación ▷ 48

Modelo conexionista

- Un conjunto de procesadores elementales densamente interconectados
- Nombres alternativos:
 - Modelo conexionista
 - Red neuronal artificial
 - Procesado distribuido y paralelo
- Neurocomputador: conjunto de procesadores hardware interconectados operando concurrentemente que implementan un modelo conexionista

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. En 1950 Rochester realiza una simulación en un computador IBM
- 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 múltiples 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 múltiples perceptrones en cascada con funciones no-lineales. El cálculo de derivadas parciales se muestra esquivo

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

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*”. Para el aprendizaje de los pesos de estas redes se recurre a una técnica de randomización conocida como “*máquina de Boltzmann restringida*” y se usa la técnica de retropropagación del error para un ajuste final. Posteriormente, en 2015 *LeCun, Bengio y Hinton* publican un artículo sobre estas técnicas en *Nature*

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

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**[†]

[†] también denominada **función logística**

Funciones de activación[†]

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

1. **LINEAL**: $g_L(z) = z$.

2. **SALTO O ESCALÓN**: $g_E(z) = \text{sgn}(z) = \begin{cases} +1 & \text{si } z > 0 \\ -1 & \text{si } z < 0 \end{cases}$

3. **RAMPA**: $g_R(z) = \begin{cases} z & \text{si } -1 < z < +1 \\ \text{sgn}(z) & \text{si } |z| \geq 1 \end{cases}$

4. **SIGMOID**: $g_S(z) = \frac{1}{1 + \exp(-z)}$

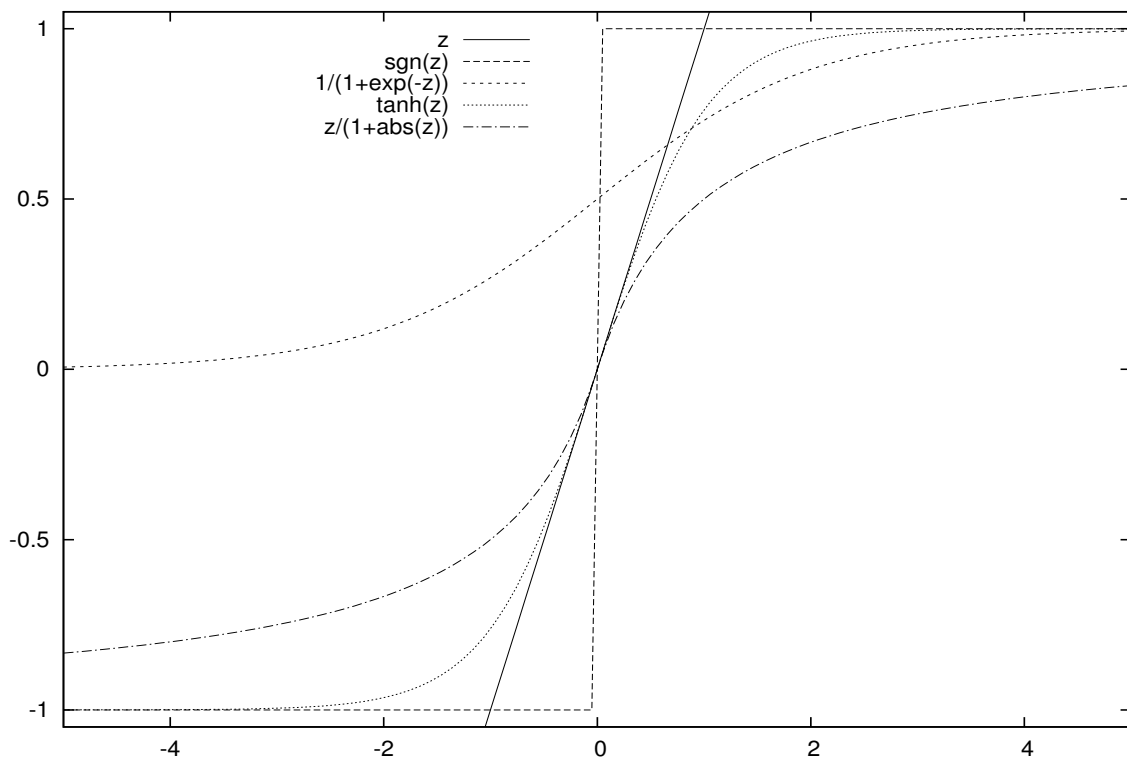
5. **TANGENTE HYPERBÓLICA**: $g_T(z) = \tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$

6. **RÁPIDA**: $g_F(z) = \frac{z}{1 + |z|}$

7. **SOFTMAX**: Dados $z_1, \dots, z_M \in \mathbb{R}$, $g_M(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^M \exp(z_j)}$ para $1 \leq i \leq M$

Ejercicio: Demostrar que $g_T(z) = 2g_S(2z) - 1 \quad \forall z \in \mathbb{R}$

Funciones de activación



Septiembre, 2016

Departament de Sistemes Informàtics i Computació

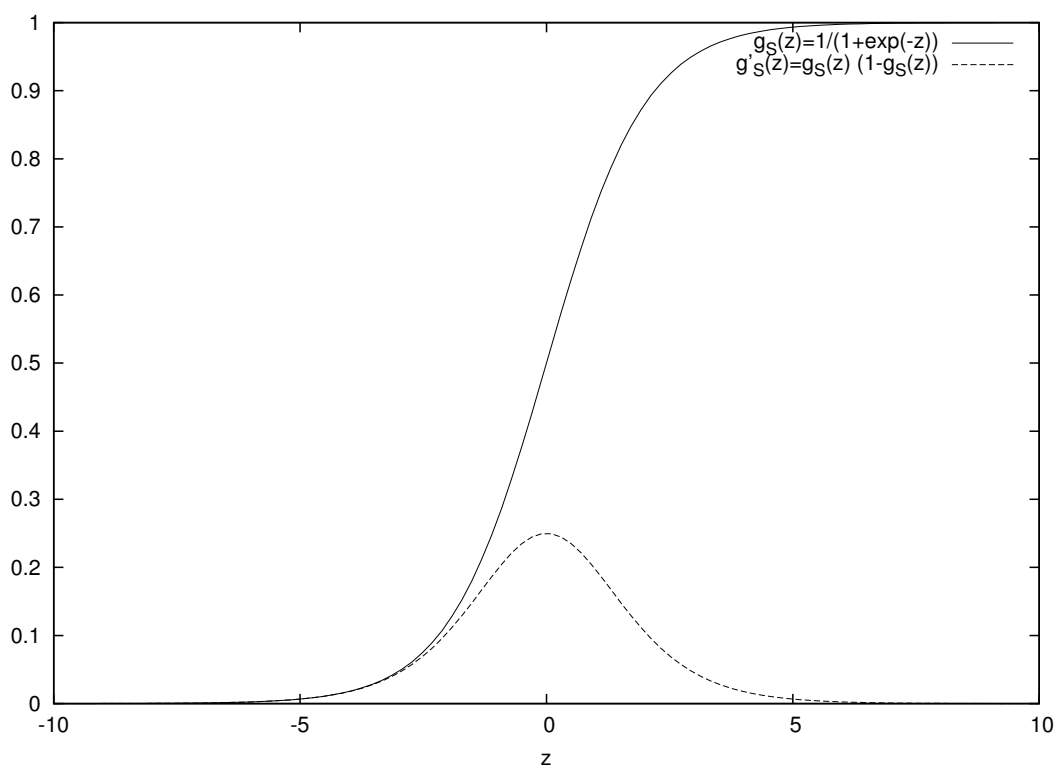
Derivadas de las funciones de activación

- **LINEAL:** $g_L(z) = z \Rightarrow g'_L(z) = \frac{d}{dz} g_L = 1$
- **ESCALÓN:** $g_E(z) = \begin{cases} 1 & \text{si } z > 0 \\ 0 & \text{si } z < 0 \end{cases} \Rightarrow g'_E(z) = \begin{cases} 0 & \text{si } z \neq 0 \\ \text{no definida} & \text{si } z = 0 \end{cases}$
- **RAMPA:** $g_R(z) = \begin{cases} z & -1 < z < +1 \\ \text{sgn}(z) & |z| \geq 1 \end{cases} \Rightarrow g'_R(z) = \begin{cases} 1 & -1 < z < +1 \\ 0 & |z| > 1 \\ \text{no def.} & |z| = 1 \end{cases}$
- **SIGMOID:** $g_S(z) = \frac{1}{1 + \exp(-z)} \Rightarrow g'_S(z) = \frac{d}{dz} g_S = g_S(z) (1 - g_S(z))$
- **TANGENTE HIPERBÓLICA:** $g_T(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)} \Rightarrow g'_T(z) = \frac{d}{dz} g_T = 1 - (g_T(z))^2$
- **RÁPIDA:** $g_F(z) = \frac{z}{1 + |z|} \Rightarrow g'_F(z) = \frac{d}{dz} g_F = \frac{1}{(1 + |z|)^2}$
- **SOFTMAX:** Para $z_1, \dots, z_n \in \mathbb{R}$, $g_M(z_i) = \frac{\exp(z_i)}{\sum_j \exp(z_j)} \Rightarrow g'_M(z_i) = g_M(z_i) (1 - g_M(z_i))$

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Derivada de la función de activación sigmoid

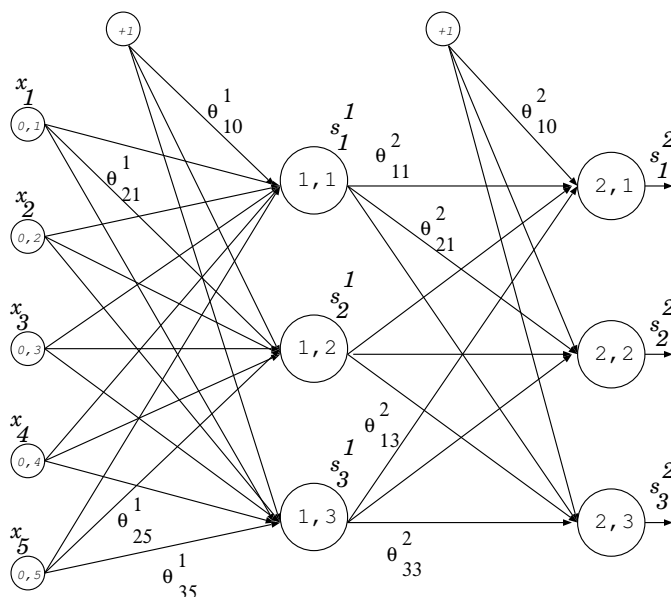


Septiembre, 2016

Departament de Sistemes Informàtics i Computació

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

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

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

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

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 Θ tal que $s^2(\mathbf{x}_n; \Theta) = \mathbf{t}_n \quad \forall n, 1 \leq n \leq N$.
- **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 $\mathbb{R}^{M_0} \rightarrow \mathbb{R}^{M_2}$:

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

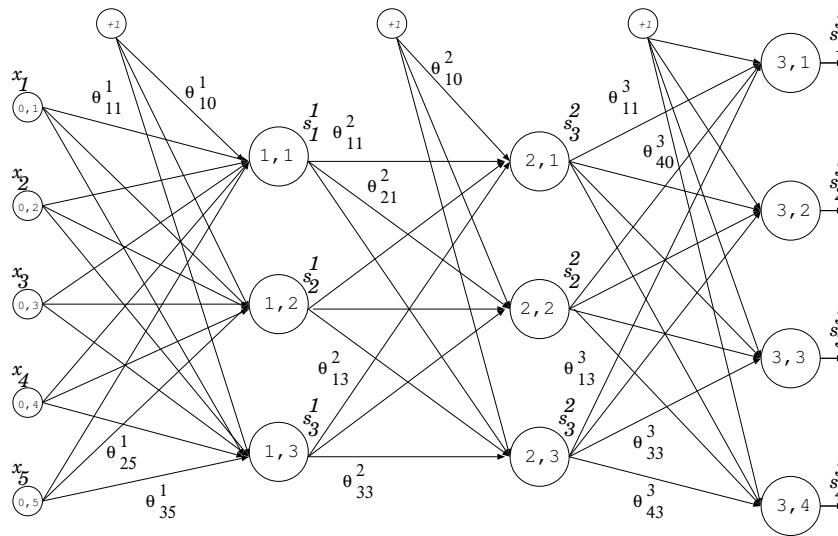
- Si todas las funciones de activación son lineales, un perceptrón multicapa define **UNA FUNCIÓN DISCRIMINANTE LINEAL**, $\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_{M_2}(\mathbf{x}))^t$:

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

- 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**:

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

Un perceptrón de tres capas

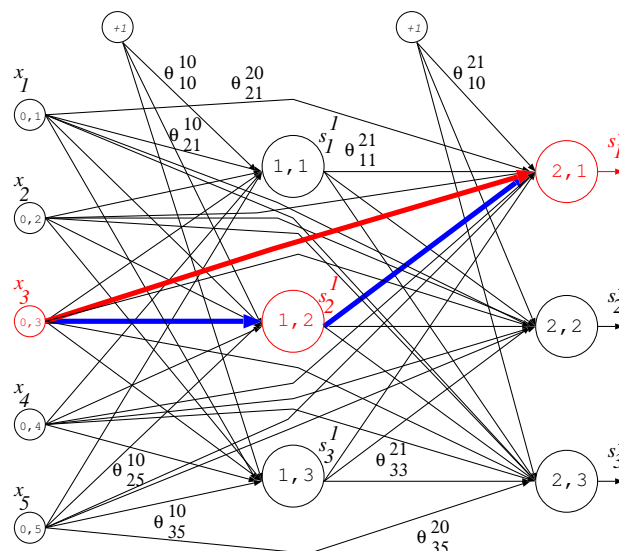


- Primera capa oculta: $s_i^1 = g(\sum_{j=0}^{M_0} \theta_{ij}^1 x_j)$ para $1 \leq i \leq M_1$
- Segunda capa oculta: $s_i^2 = g(\sum_{j=0}^{M_1} \theta_{ij}^2 s_j^1)$ para $1 \leq i \leq M_2$
- Capa de salida: $s_i^3 = g(\sum_{j=0}^{M_2} \theta_{ij}^3 s_j^2)$ para $1 \leq i \leq M_3$

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

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

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

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 $d' \equiv M_2$)

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

- 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*

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

El perceptrón multicapa como clasificador

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

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

- 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]

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Index

- 1 Redes neuronales multicapa ▷ 2
- 2 *Algoritmo de retropropagación del error (BackProp)* ▷ 18
- 3 Variantes del BackProp ▷ 35
- 4 Redes neuronales radiales ▷ 39
- 5 Aplicaciones ▷ 43
- 6 Notación ▷ 48

Aprendizaje de los pesos de un perceptrón multicapa

PROBLEMA: Dada la topología de un perceptrón multicapa y 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 $\Theta \in \mathbb{R}^D : s^2(\mathbf{x}_n; \Theta) = \mathbf{t}_n$ para $1 \leq n \leq N$ (o $s^2(\mathbf{x}_n; \Theta) \approx \mathbf{t}_n$).

- Resolver el sistema de ecuaciones no lineales $s^2(\mathbf{x}_n) = \mathbf{t}_n, 1 \leq n \leq N$: raramente tiene solución.
- **Problema reformulado:** Dada la topología de un perceptrón multicapa y 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 Θ que minimice el error cuadrático medio:

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

- Solución: **DESCENSO POR GRADIENTE:**

$$\Delta \theta_{ij}^l = -\rho \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^l}$$

$$1 \leq l \leq 2, \quad 1 \leq i \leq M_l, \quad 0 \leq j \leq M_{l-1}$$

Algoritmo de retropropagación del error (BackProp)

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

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

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

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

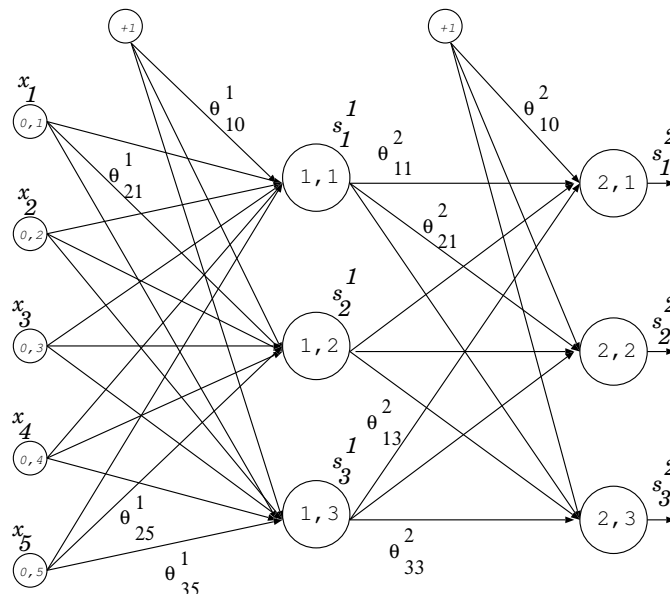
$$\Delta\theta_{ij}^1 = -\rho \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^1} = \rho \sum_{n=1}^N \delta_i^1(\mathbf{x}_n) x_{nj}$$

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

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Algoritmo BackProp: cálculo hacia adelante



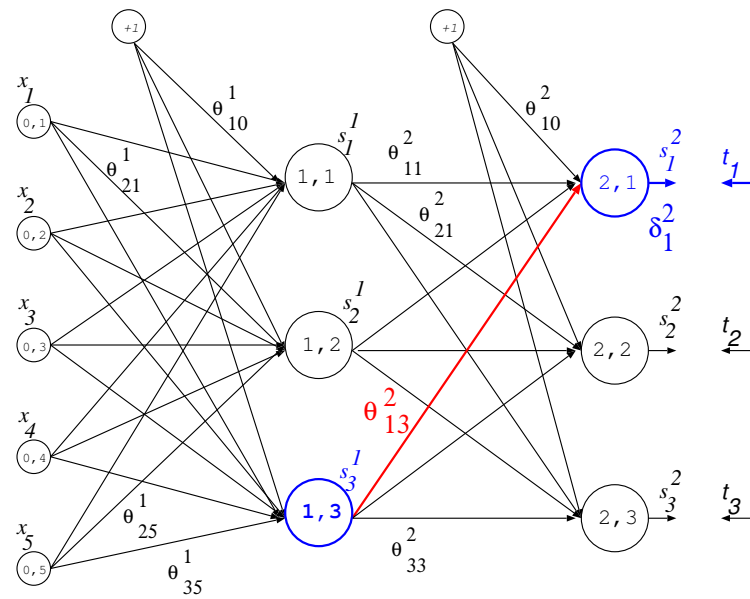
$$s_i^1(\mathbf{x}) = g(\phi_i^1) = g\left(\sum_{j=0}^{M_0} \theta_{ij}^1 x_j\right) \quad 1 \leq i \leq M_1; \quad s_j^2 = g(\phi_j^2) = g\left(\sum_{k=0}^{M_1} \theta_{jk}^2 s_k^1(\mathbf{x})\right), \quad 1 \leq j \leq M_2$$

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

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

BackProp: Actualización pesos de la capa de salida



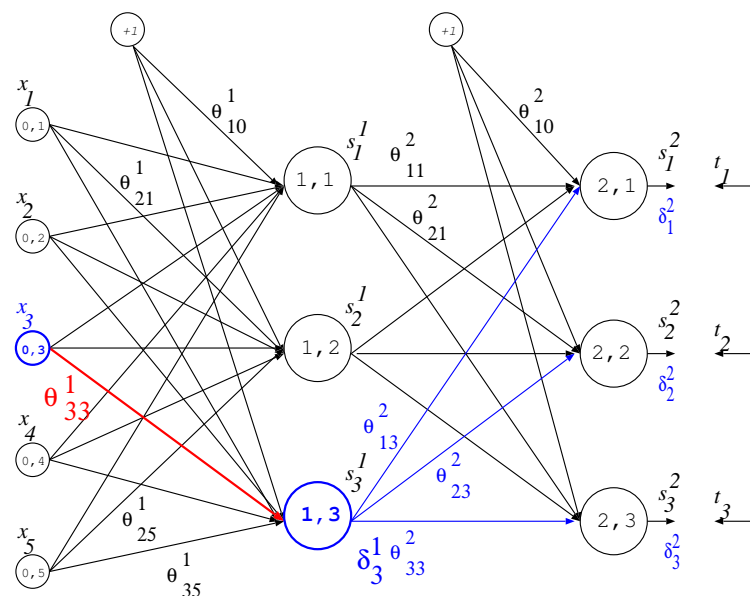
$$\Delta \theta_{13}^2 = \rho \delta_1^2 s_3^1 = \rho (t_1 - s_1^2) g'(\phi_1^2) s_3^1$$

(para solo una muestra de entrenamiento)

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

BackProp: Actualización pesos de la capa oculta



$$\Delta \theta_{33}^1 = \rho \delta_3^2 x_3 = \rho \left(\sum_{r=1}^{M_2} \delta_r^2 \theta_{r3}^2 \right) g'(\phi_3^1) x_3$$

(para solo una muestra de entrenamiento)

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Regla de la cadena en el cálculo de derivadas

- Función simple de otra función de una variable: $f, g : \mathbb{R} \rightarrow \mathbb{R}$:

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

- Función de dos funciones de M variables: $g_1, g_2 : \mathbb{R}^M \rightarrow \mathbb{R}, f : \mathbb{R}^2 \rightarrow \mathbb{R}$

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

- Función de N funciones de M variables: $g_1, \dots, g_N : \mathbb{R}^M \rightarrow \mathbb{R}, f : \mathbb{R}^N \rightarrow \mathbb{R}$:

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

Derivación del algoritmo BackProp (I)

Actualización de los pesos de la capa de salida θ_{kl}^2 (una sola muestra, $N = 1$)

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

$$\frac{\partial q_S(\Theta)}{\partial \theta_{ij}^2} = \frac{\partial q_S(\Theta)}{\partial \phi_i^2} \frac{\partial \phi_i^2}{\partial \theta_{ij}^2} = \frac{\partial q_S(\Theta)}{\partial s_i^2} \frac{\partial s_i^2}{\partial \phi_i^2} \frac{\partial \phi_i^2}{\partial \theta_{ij}^2} = [(-1)(t_i - s_i^2)] g'(\phi_i^2) s_j^1$$

$$\delta_i^2 \stackrel{\text{def}}{=} (t_i - s_i^2) g'(\phi_i^2) \left(= -\frac{\partial q_S(\Theta)}{\partial \phi_i^2} \right) \Rightarrow \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^2} = -\delta_i^2 s_j^1$$

$$\Delta \theta_{ij}^2 = -\rho \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^2} = \rho \delta_i^2 s_j^1 \quad 1 \leq i \leq M_2, 0 \leq j \leq M_1$$

Derivación del algoritmo BackProp (II)

Actualización de los pesos de la capa de la capa oculta θ_{ij}^1 (para $N = 1$)

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

$$\frac{\partial q_S(\Theta)}{\partial \theta_{ij}^1} = \frac{\partial q_S(\Theta)}{\partial \phi_i^1} \frac{\partial \phi_i^1}{\partial \theta_{ij}^1} = \frac{\partial q_S(\Theta)}{\partial s_i^1} \frac{\partial s_i^1}{\partial \phi_i^1} \frac{\partial \phi_i^1}{\partial \theta_{ij}^1}$$

$$= \left(\sum_{l=0}^{M_2} \frac{\partial q_S(\Theta)}{\partial \phi_l^2} \frac{\partial \phi_l^2}{\partial s_i^1} \right) g'(\phi_i^1) x_j = \left(\sum_{l=0}^{M_2} -\delta_l^2 \theta_{li}^2 \right) g'(\phi_i^1) x_j$$

$$\delta_i^1 \stackrel{\text{def}}{=} \left(\sum_{l=1}^{M_2} \delta_l^2 \theta_{li}^2 \right) g'(\phi_i^1) \implies \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^1} = -\delta_i^1 x_j$$

$$\Delta \theta_{ij}^1 = -\rho \frac{\partial q_S(\Theta)}{\partial \theta_{ij}^1} = \rho \delta_i^1 x_j \quad 1 \leq i \leq M_1, 0 \leq j \leq M_0$$

BackProp para perceptrones de tres capas

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

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

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

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

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

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

Algoritmo BACKPROP

Entrada: Una topología, datos de entrenamiento S , un factor de aprendizaje ρ , pesos iniciales θ_{ij}^l $1 \leq l \leq L$, $1 \leq i \leq M_l$, $0 \leq j \leq M_{l-1}$ y condiciones de convergencia

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

Método:

Mientras no converja

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 toda muestra de aprendizaje $(\mathbf{x}_n, \mathbf{t}_n) \in S$

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

Calcular $\phi_i^l(\mathbf{x}_n)$ y $s_i^l(\mathbf{x}_n) = g(\phi_i^l(\mathbf{x}_n))$ para $1 \leq i \leq M_l$

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

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

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

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

Fin para

Fin para

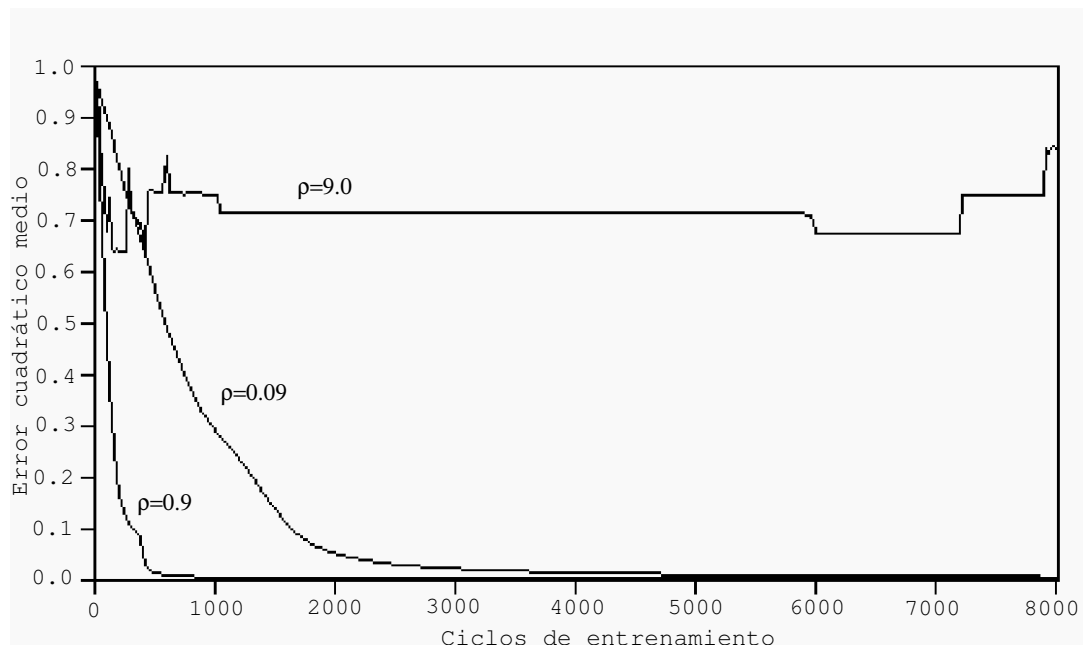
Actualizar pesos: $\theta_{ij}^l = \theta_{ij}^l + \rho \Delta\theta_{ij}^l$

Fin mientras

El coste computacional del algoritmo BackProp es $O(N D)$ en cada iteración.

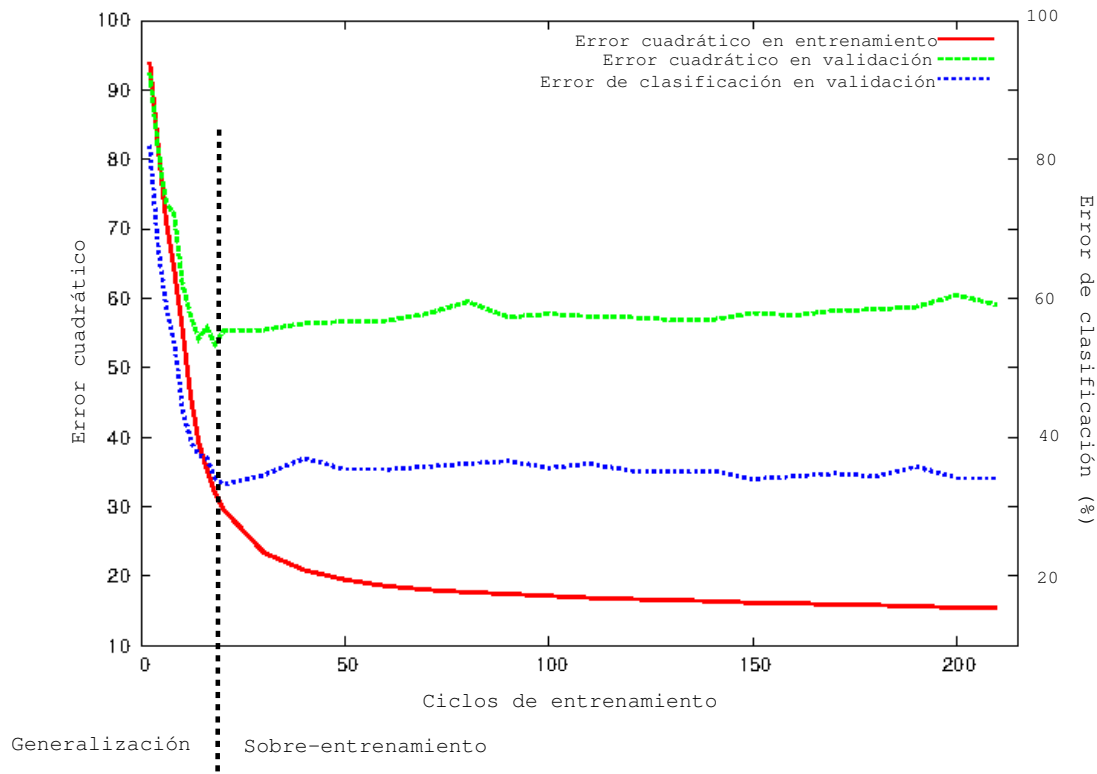
DEMO

La selección del factor de aprendizaje



DEMO

Condiciones de convergencia



Codificación

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

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

- 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
- **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.*

En ocasiones, se suele utilizar las salidas del perceptrón como una aproximación a esa distribución a-posteriori (p.e. para un perceptrón de una capa oculta): $s_i^2(\mathbf{x}; \Theta) \approx P(i | \mathbf{x})$ para lo que se suele utilizar la función “softmax” como función de activación en la capa de salida, lo que permite su normalización.

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Estimación de la topología

- **MÉTODOS DE PODA:**
 - Poda de pesos:
 - * basada en **RELEVANCIA**
 - * basada en **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

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Aspectos computacionales del aprendizaje con el PM

- Sobre la intratabilidad del aprendizaje de redes neuronales (Blum and Rivest, 1992): La búsqueda de un conjunto de pesos (rationales) tal que el perceptrón multicapa (con funciones de activación en escalón) genere la salida deseada (en el caso de que sea posible) para cada muestra de entrenamiento es **NP-duro**.
- Sobre el tamaño del conjunto de entrenamiento (Ripley, 1993 y Tema 2.1): Si el perceptrón multicapa tiene M nodos y un número arbitrario de capas, el número de muestras de entrenamiento debe ser proporcional a $\frac{D}{\epsilon} \log \frac{M}{\epsilon}$ donde ϵ 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 *Variantes del BackProp* ▷ 35
- 4 Redes neuronales radiales ▷ 39
- 5 Aplicaciones ▷ 43
- 6 Notación ▷ 48

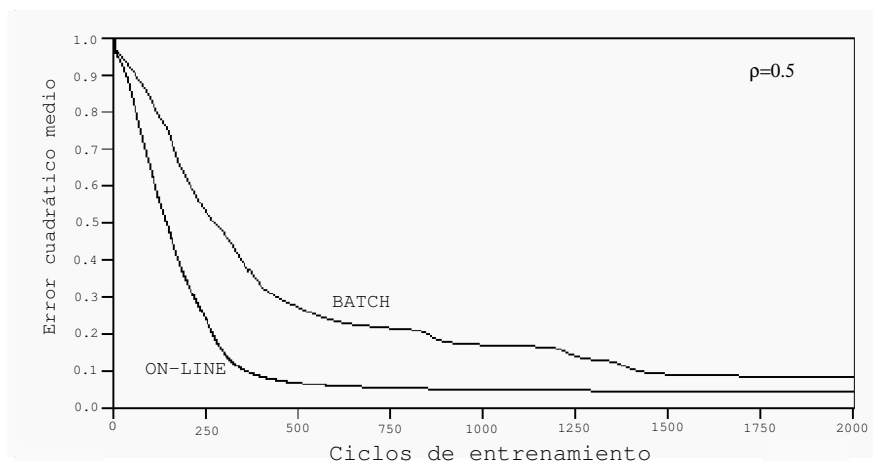
BackProp incremental u “on-line”

- BackProp “batch” (en cada iteración k se procesan todas las $\mathbf{x}_n, 1 \leq n \leq N$):

$$\Delta\theta_{ij}^l(k) = \rho_k \sum_{n=1}^N \delta_i^l(\mathbf{x}_n) s_j^{l-1}(\mathbf{x}_n) \quad 1 \leq l \leq L, 1 \leq i \leq M_l, 1 \leq j \leq M_{l-1}$$

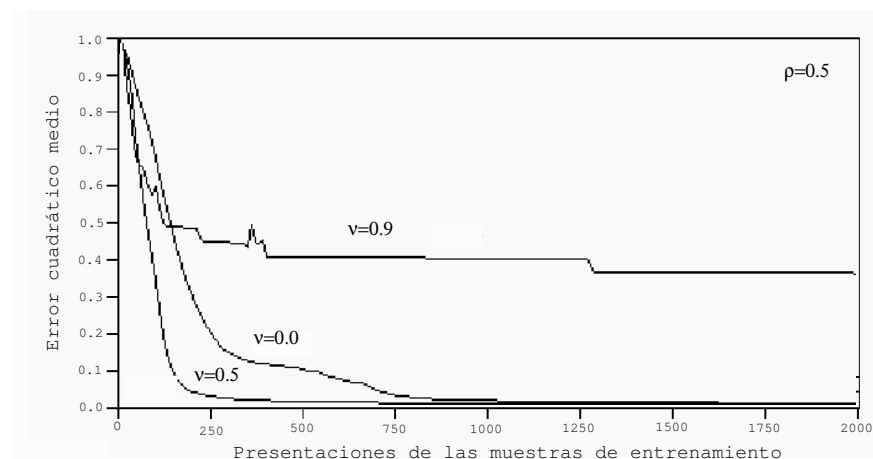
- BackProp “on-line” (en la iteración k -ésima se procesa solo $\mathbf{x}_{k'}, k' = 1 + k \bmod N$):

$$\Delta\theta_{ij}^l(k) = \rho_k \delta_i^l(\mathbf{x}_{k'}) s_j^{l-1}(\mathbf{x}_{k'}) \quad 1 \leq l \leq L, 1 \leq i \leq M_l, 1 \leq j \leq M_{l-1}$$



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

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

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

$$\Delta \theta_{ij}^l(k) = \rho_k \sum_{n=1}^N \delta_i^l(\mathbf{x}_n) s_j^{l-1}(\mathbf{x}_n) + 2 \rho_k \lambda \theta_{ij}^l(k)$$

Redes neuronales radiales

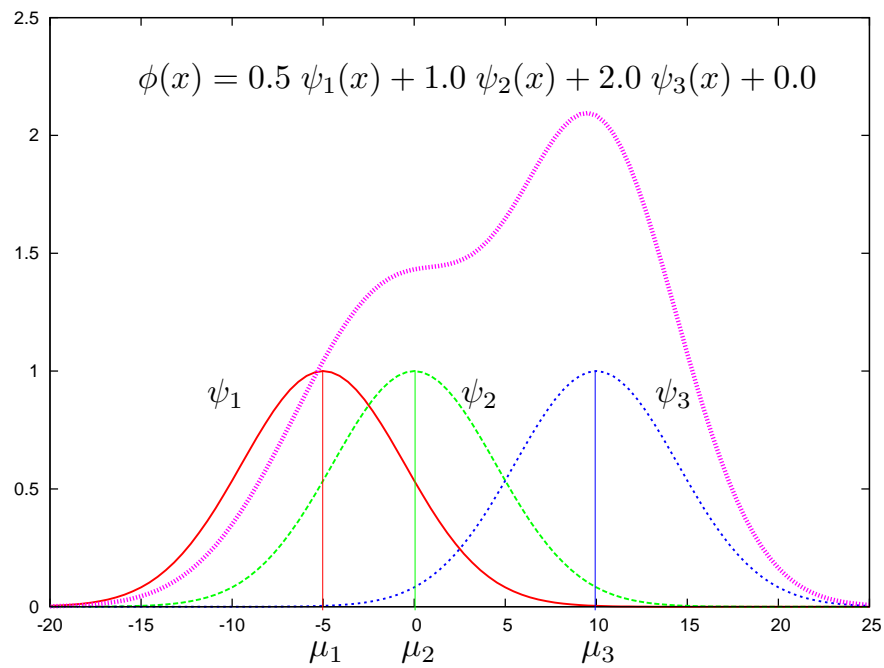
- **Una red neuronal radial** (RNR) es una FDLG de la forma:

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

donde: $\mathbf{x} \in \mathbb{R}^d$, $\boldsymbol{\mu}_i \in \mathbb{R}^d$, para $1 \leq i \leq d'$ $\theta_m \in \mathbb{R}^{d'}$, $\theta_{m0} \in \mathbb{R}$ y ψ 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).

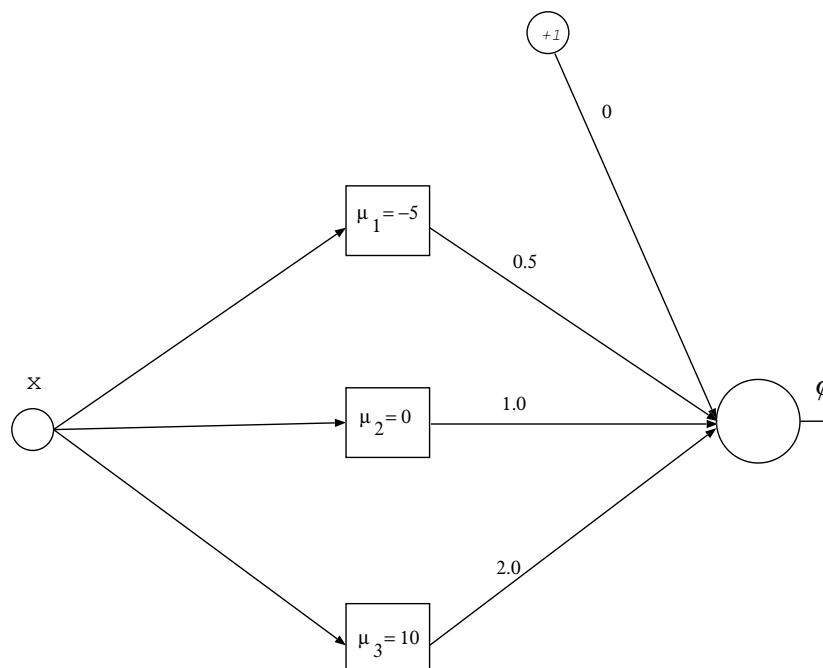
Ejemplo de funciones básicas radiales en \mathbb{R}



Septiembre, 2016

Departament de Sistemes Informàtics i Computació

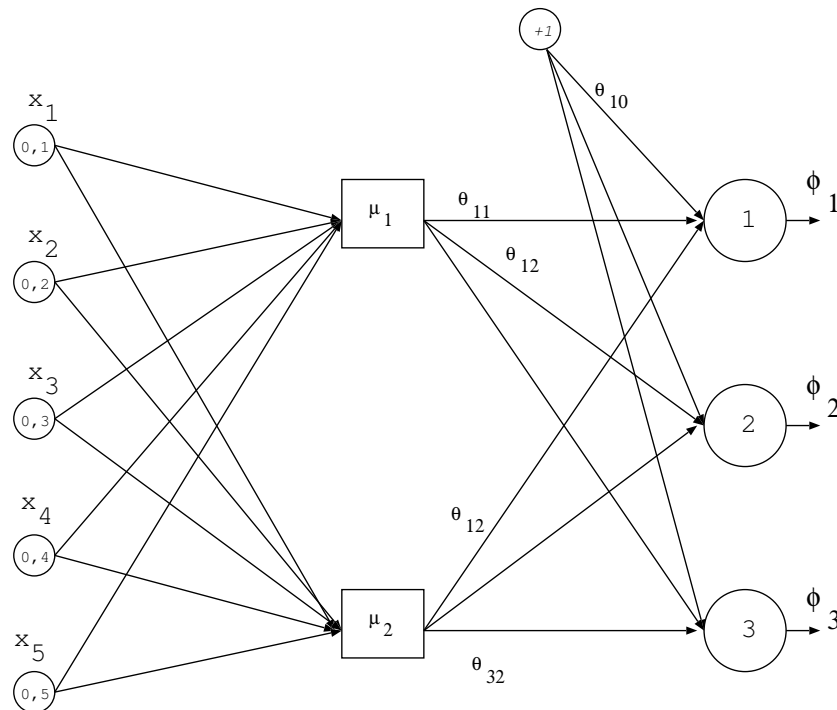
Ejemplo de funciones básicas radiales en \mathbb{R}



Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Redes neuronales radiales



Index

- 1 Redes neuronales multicapa ▷ 2
- 2 Algoritmo de retropropagación del error (BackProp) ▷ 18
- 3 Variantes del BackProp ▷ 35
- 4 Redes neuronales radiales ▷ 39
- 5 *Aplicaciones* ▷ 43
- 6 Notación ▷ 48

Ejemplos

Clasificación de dígitos I

Clasificación de dígitos II

Predicción

Compresión de imágenes

<http://neuron.eng.wayne.edu/bpImageCompression9PLUS/bp9PLUS.html>

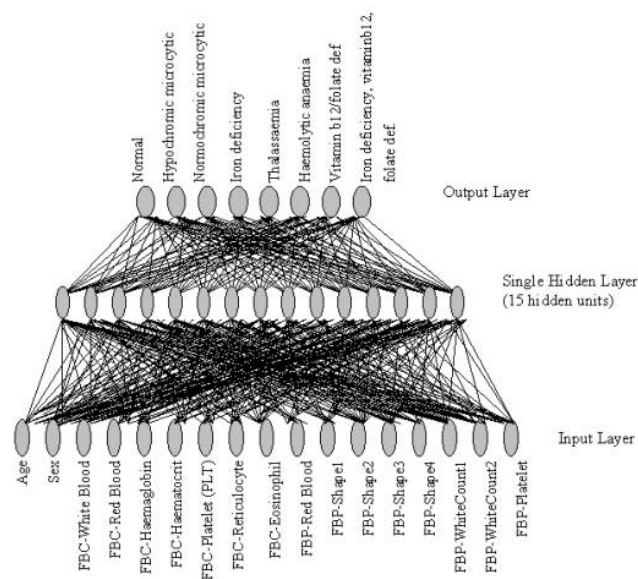
Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Aprendizaje Automático. 2016-2017

Redes Neuronales Multicapa: 6.45

Ejemplo: Clasificación de anemias



<http://www.generation5.org/content/2004/NNinAnaemia.asp>

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

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)

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Algunas aplicaciones

- Aplicaciones varias
<http://www.calsci.com/Applications.html>
- Diagnóstico médico
<http://www.generation5.org/content/2004/MedicalDiagnosis.asp>
- Anemias
<http://www.generation5.org/content/2004/NNinAnaemia.asp>
- Aplicaciones en física de altas energías
<http://neuralnets.web.cern.ch/NeuralNets/nnwInHep.html>
- Demos en java
<http://lcn.epfl.ch/tutorial/english/index.html>
<http://www.aissance.org/downloads.shtml>

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Notación

- **Funciones discriminantes lineales:** $\phi(\mathbf{x}; \Theta) = \Theta^t \mathbf{x}$ para una entrada \mathbf{x} y parámetros Θ compuestos por vector de pesos y umbral (θ, θ_0)
- **Funciones discriminantes lineales con activación:** $g \circ \phi(\mathbf{x}; \theta)$ para una entrada \mathbf{x} , parámetros (θ, θ_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: θ_{ij}^k . Θ es un vector de talla D formado por todos los pesos θ_{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:** ρ, ν y λ
- **Media y desviación típica:** μ y σ