

5.4 Redes neuronales para clasificación

Inteligencia Artificial

3er curso INF



Aprendizaje supervisado: clasificación

Aprendizaje automático por **inducción a partir de datos etiquetados**

$$\mathcal{D} = \{(\mathbf{x}_n, c_n)\}_{n=1}^N$$

$\mathbf{x} \in \mathcal{X}$: vector de **atributos**, características, variables independientes, variables de entrada, covariables ...

Término bias
 $x_{n0} = 1$

$$\mathbf{x}_n^T = (x_{n0} \ x_{n1} \ x_{n2} \ \dots \ x_{nD}) \ [(D + 1)\text{-vector dimensional}]$$

$c \in \mathcal{C}$: (clase) **etiqueta**, dependent variable, outcome, target, ...

Clasificación: Etiquetas c discretas (e.g. $c \in \{C_1, \dots, C_K\}$)

Algoritmo de aprendizaje : \mathcal{L}

\mathbf{w} : parámetros del modelo

$$\mathcal{L}: \mathcal{D} = \{(\mathbf{x}_n, c_n)\}_{n=1}^N \rightarrow h(\cdot; \mathbf{w})$$

$$h(\cdot; \mathbf{w}): \mathbf{x} \in \mathcal{X} \rightarrow h(\mathbf{x}; \mathbf{w}) \in \{C_1, \dots, C_K\}$$

Clasificación binaria

n	x_{n1}	x_{n2}	...	x_{nD}	C_n
1	2.3	0	...	10.3	C_0
2	2.5	1	...	13.1	C_1
3	2.6	0	...	-2.7	C_1
4	2.7	-1	...	-5.4	C_0
5	2.9	0	...	2.1	C_1
6	3.1	0	...	-10.9	C_0

Clasificación binaria

Codificación 0-1 de los **datos etiquetados**

$$\mathcal{D} = \{(\mathbf{x}_n, c_n)\}_{n=1}^N \Rightarrow \mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^N; \quad t_n = \begin{cases} 0 & \text{si } c_n = C_0 \\ 1 & \text{si } c_n = C_1 \end{cases}$$

$$\mathcal{L}: \mathcal{D} = \{(\mathbf{x}_n, c_n)\}_{n=1}^N \rightarrow o(\cdot; \mathbf{w})$$

$$z(\cdot; \mathbf{w}): \mathbf{x} \in \mathcal{X} \rightarrow z(\mathbf{x}; \mathbf{w}) \in \mathbb{R}$$

$$o(\mathbf{x}; \mathbf{w}) = \varphi^{(o)}(z(\mathbf{x}; \mathbf{w})) \in [0,1]$$

- La salida $o(\mathbf{x}; \mathbf{w})$ es una estimación de la clase C_1 posterior

Modelo discriminativo $\hat{p}(C_1 | \mathbf{x}, \mathbf{w}) = \varphi^{(o)}(z(\mathbf{x}; \mathbf{w}))$

Regresión logística : $\varphi^{(o)}(z) = \sigma(z) = \frac{1}{1+e^{-z}}$

Regresión probit: $\varphi^{(o)}(z) = \Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{y^2}{2}} dy$

Clasificación binaria: codificación 0-1

n	x_{n1}	x_{n2}	...	x_{nD}	t_n
1	2.3	0	...	10.3	0
2	2.5	1	...	13.1	1
3	2.6	0	...	-2.7	1
4	2.7	-1	...	-5.4	0
5	2.9	0	...	2.1	1
6	3.1	0	...	-10.9	0

Aprendizaje por máxima probabilidad

Se suponen
Muestras iid
(distribuidas
Idéntica e
Independiente-
mente)

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \mathcal{L}(\mathbf{w})$$

- Probabilidad del modelo, dado $\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^N$

$$\mathcal{L}(\mathbf{w}) = \hat{P}(\{t_n\}_{n=1}^N | \{\mathbf{x}_n\}_{n=1}^N, \mathbf{w}) = \prod_{n=1}^N \hat{p}(t_n | \mathbf{x}_n, \mathbf{w}) =$$

Factoriza porque
se asumen muestras
independientes

$$= \prod_{n=1}^N (\hat{p}(C_0 | \mathbf{x}_n, \mathbf{w}))^{1-t_n} (\hat{p}(C_1 | \mathbf{x}_n, \mathbf{w}))^{t_n}$$

Estimación de la
clase posterior

Misma distribución:
se asumen muestras
idénticamente dist.

- Función de verosimilitud

$$\log \mathcal{L}(\mathbf{w}) = \sum_{n=1}^N ((1 - t_n) \log \hat{p}(C_0 | \mathbf{x}_n, \mathbf{w}) + t_n \log \hat{p}(C_1 | \mathbf{x}_n, \mathbf{w}))$$

Mismo maximizador que $\mathcal{L}(\mathbf{w})$ (log is monotono)

Minimización del error de entropía cruzada

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} CE(\mathbf{w})$$

- Error de entropía cruzada

- $CE(\mathbf{w}) = -\log \mathcal{L}(\mathbf{w})$

Mismo maximizador que $\mathcal{L}(\mathbf{w})$

$$= -\sum_{n=1}^N \left((1 - t_n) \log \hat{p}(C_0 | \mathbf{x}_n, \mathbf{w}) + t_n \log \hat{p}(C_1 | \mathbf{x}_n, \mathbf{w}) \right)$$

$$= -\sum_{n=1}^N \left((1 - t_n) \log(1 - \hat{p}(C_1 | \mathbf{x}_n, \mathbf{w})) + t_n \log \hat{p}(C_1 | \mathbf{x}_n, \mathbf{w}) \right)$$

- Modelo discriminativo : $\hat{p}(C_1 | \mathbf{x}, \mathbf{w}) = \varphi^{(o)}(z(\mathbf{x}_n; \mathbf{w}))$

$$CE(\mathbf{w}) = -\sum_{n=1}^N \left((1 - t_n) \log \left(1 - \varphi^{(o)}(z(\mathbf{x}_n; \mathbf{w})) \right) + t_n \log \left(\varphi^{(o)}(z(\mathbf{x}_n; \mathbf{w})) \right) \right)$$

$\hat{p}(C_1 | \mathbf{x}, \mathbf{w})$

$\hat{p}(C_0 | \mathbf{x}, \mathbf{w})$

Perceptrón monocapa

Puede resolver solo problemas linealmente separables

$$\mathbf{w}^T = (w_0, w_1, \dots, w_D)$$
$$\mathbf{x}^T = (x_0 \ x_1 \ \dots \ x_D)$$
$$\mathbf{w}^T \mathbf{x} = \sum_{d=0}^D w_d x_d$$

Capa de entrada

Término bias

$x_0 = 1$

w_0

x_1

w_1

x_2

w_2

\vdots

w_D

x_D

Entradas

Capa de salida

$\Sigma \varphi^{(o)}$

$$o(\mathbf{x}; \mathbf{w}) = \varphi^{(o)}(\mathbf{w}^T \mathbf{x})$$

La salida es una estimación de $p(C_1|\mathbf{x})$

Aprendizaje: Se determinan los pesos de la red por minimización de una función de coste. Ej:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left(- \sum_{n \in C_0} \log(1 - o(\mathbf{x}_n; \mathbf{w})) - \sum_{n \in C_1} \log(o(\mathbf{x}_n; \mathbf{w})) \right)$$

Error de entropía cruzada

Regresión logística

- Clase posteriores

$$p(C_1|\mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{-z(\mathbf{x}; \mathbf{w})}}$$

$$p(C_0|\mathbf{x}, \mathbf{w}) = 1 - p(C_1|\mathbf{x}, \mathbf{w}) = \frac{e^{-z(\mathbf{x}; \mathbf{w})}}{1 + e^{-z(\mathbf{x}; \mathbf{w})}} = \frac{1}{1 + e^{z(\mathbf{x}; \mathbf{w})}}$$

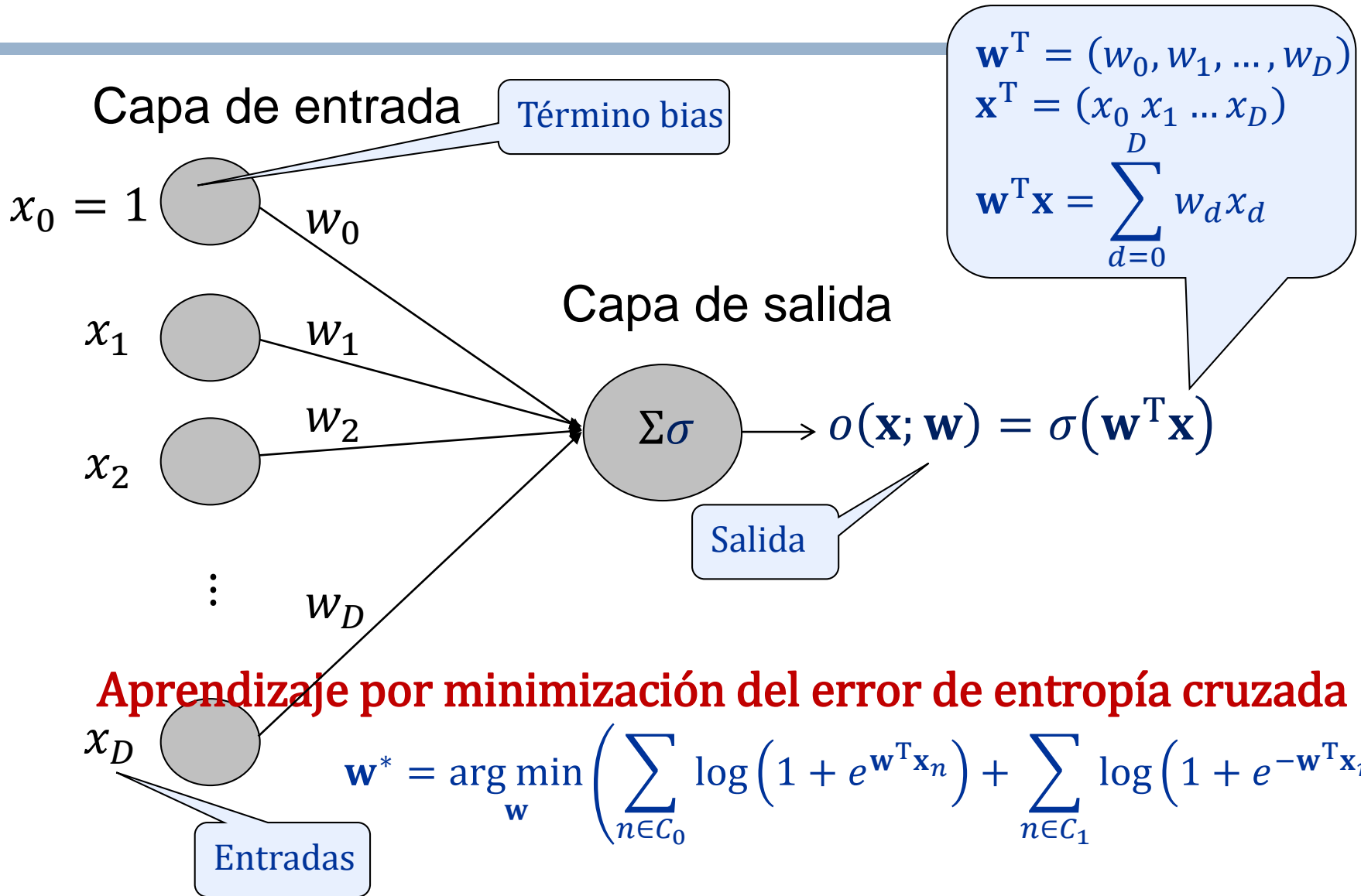
- Probabilidades-log (log-odds): $z(\mathbf{x}; \mathbf{w}) = \log \frac{p(C_1|\mathbf{x}, \mathbf{w})}{p(C_0|\mathbf{x}, \mathbf{w})}$

- Modelo lineal para log-odds: $z(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} = \sum_{d=0}^D w_d x_d$

$$\begin{aligned} \mathbf{w}^T &= (w_0, w_1, \dots, w_D) \\ \mathbf{x}^T &= (x_0 \ x_1 \ \dots \ x_D) \end{aligned}$$

Perceptrón monocapa!

Perceptrón monocapa: regresión logística



Error de entropía cruzada para regresión logística

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} CE(\mathbf{w})$$

$$\frac{\partial \sigma(\mathbf{w}^T \mathbf{x}_n)}{\partial \mathbf{w}} = \mathbf{x}_n \sigma(\mathbf{w}^T \mathbf{x}_n) (1 - \sigma(\mathbf{w}^T \mathbf{x}_n))$$

■ Error de entropía cruzada

$$CE(\mathbf{w}) = - \sum_{n=1}^N \left((1 - t_n) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_n)) + t_n \log \sigma(\mathbf{w}^T \mathbf{x}_n) \right)$$

■ Gradiente del error de entropía cruzada

$$\frac{\partial}{\partial \mathbf{w}} CE(\mathbf{w}) = \sum_{n=1}^N \left((1 - t_n) \mathbf{x}_n \sigma(\mathbf{w}^T \mathbf{x}_n) - t_n \mathbf{x}_n (1 - \sigma(\mathbf{w}^T \mathbf{x}_n)) \right)$$

$$\delta_n = \sigma(\mathbf{w}^T \mathbf{x}_n) - t_n$$

$$= \sum_{n=1}^N (\sigma(\mathbf{w}^T \mathbf{x}_n) - t_n) \mathbf{x}_n = \sum_{n=1}^N \delta_n \mathbf{x}_n$$

Perceptrón de una capa: aprendizaje por lotes

Los atributos deberían
Estar escalados!

ENTRADA: Instancias de entrenamiento : $\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^N$

Parámetro de aprendizaje : $\eta > 0$

SALIDA: $\mathbf{w}^* = \arg \min_{\mathbf{w}} CE(\mathbf{w})$

Similar a Rosenblatt!
 $\delta_n = \sigma(\mathbf{w}^T \mathbf{x}_n) - t_n$
predicción error

1. Inicializar aleatoriamente $\mathbf{w} \sim U[-0.5, 0.5]^{(D+1)}$

2. $n_{epoch} = 0$

3. Mientras no se cumplan los criterios de convergencia

3.1 Incrementa contador de épocas: $n_{epoch} = n_{epoch} + 1$

3.2 Calcula las salidas de red : $\sigma(\mathbf{w}^T \mathbf{x}_n); \quad n = 1, \dots, N$

3.3 Calcula el gradiente : $\frac{\partial}{\partial \mathbf{w}} CE(\mathbf{w}) = \sum_{n=1}^N \delta_n \mathbf{x}_n$

3.4 Actualizar pesos : $\mathbf{w} = \mathbf{w} - \eta \sum_{n=1}^N \delta_n \mathbf{x}_n$

Perceptrón de una capa: aprendizaje en línea

ENTRADA: Instancias de entrenamiento : $\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^N$

Parámetro de aprendizaje : $\eta > 0$

SALIDA: $\mathbf{w}^* = \arg \min_{\mathbf{w}} CE(\mathbf{w})$

1. Inicializa aleatoriamente $\mathbf{w} \sim U[-0.5, 0.5]^{(D+1)}$

2. $n_{epoch} = 0$

3. Mientras no se cumplan los criterios de convergencia

3.1 Incrementa contador de épocas : $n_{epoch} = n_{epoch} + 1$

3.2 Para $n = 1, \dots, N$

Calcula la salida de red : $\sigma(\mathbf{w}^T \mathbf{x}_n)$

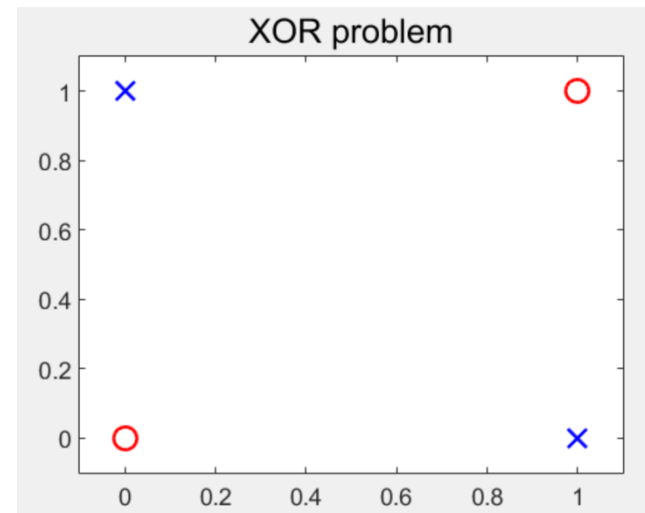
Calcula error de predicción : $\delta_n = \sigma(\mathbf{w}^T \mathbf{x}_n) - t_n$

Actualiza pesos: $\mathbf{w} = \mathbf{w} - \eta \delta_n \mathbf{x}_n$

XOR: Un problema no separable linealmente

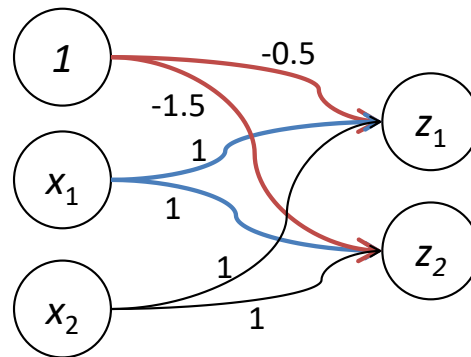
- El perceptrón de una sola capa solo puede abordar problemas que son linealmente separables
- Por lo tanto, el simple problema XOR, que no es linealmente separable, no se puede resolver con esta máquina de aprendizaje.

x_1	x_2	h
0	0	1
0	1	0
1	0	0
1	1	1

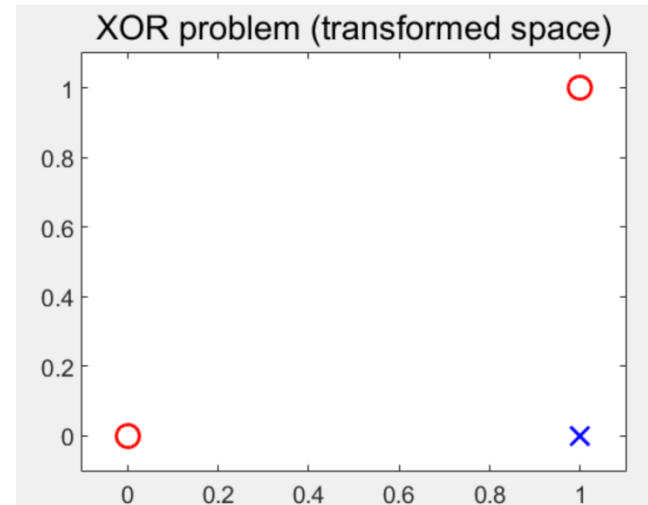


XOR: Construcción de características no lineales

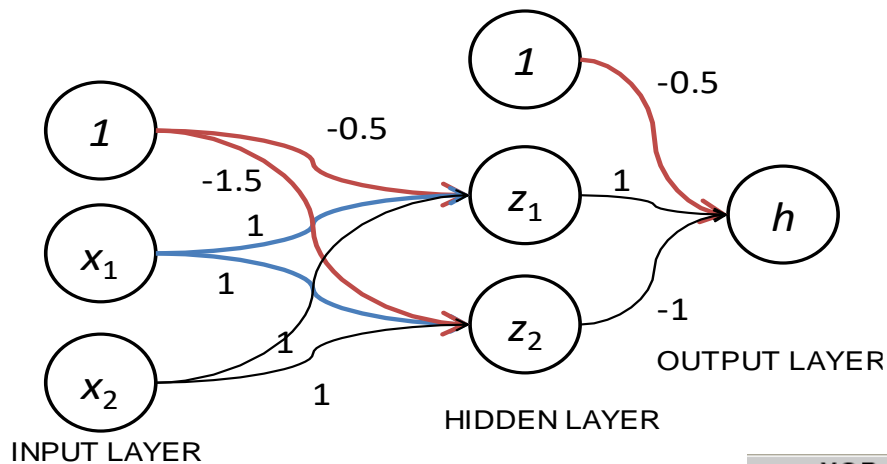
- Considere el problema XOR en un espacio de características transformado



x_1	x_2	z_1	z_2	h
0	0	0	0	0
0	1	1	0	1
1	0	1	0	1
1	1	1	1	0



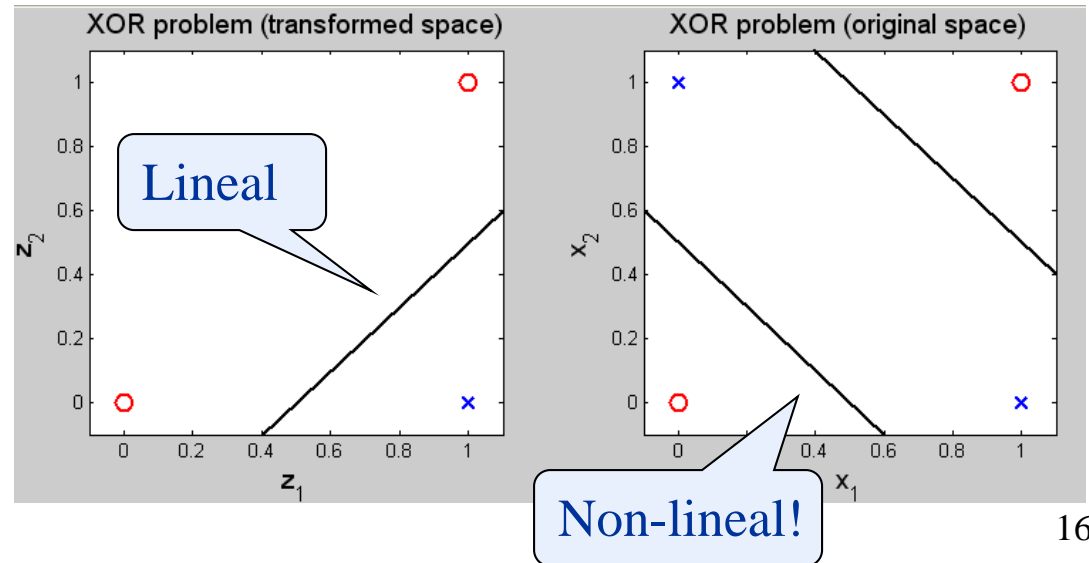
XOR: Un modelo lineal en espacio de características



x_1	x_2	z_1	z_2	h
0	0	0	0	0
0	1	1	0	1
1	0	1	0	1
1	1	1	1	0

$$\begin{cases} z_1 = \theta(x_1 + x_2 - 0.5) \\ z_2 = \theta(x_1 + x_2 - 1.5) \end{cases}$$

$$h = \theta(z_1 - z_2 - 0.5)$$



Perceptrón multicapa

$$\mathbf{x} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_D \end{pmatrix}$$

vector de atributos
(de entrada)

Oculto a la entrada
pesos de capa

$$\mathbf{V} = (\mathbf{v}_1 \mathbf{v}_2 \dots \mathbf{v}_J)$$

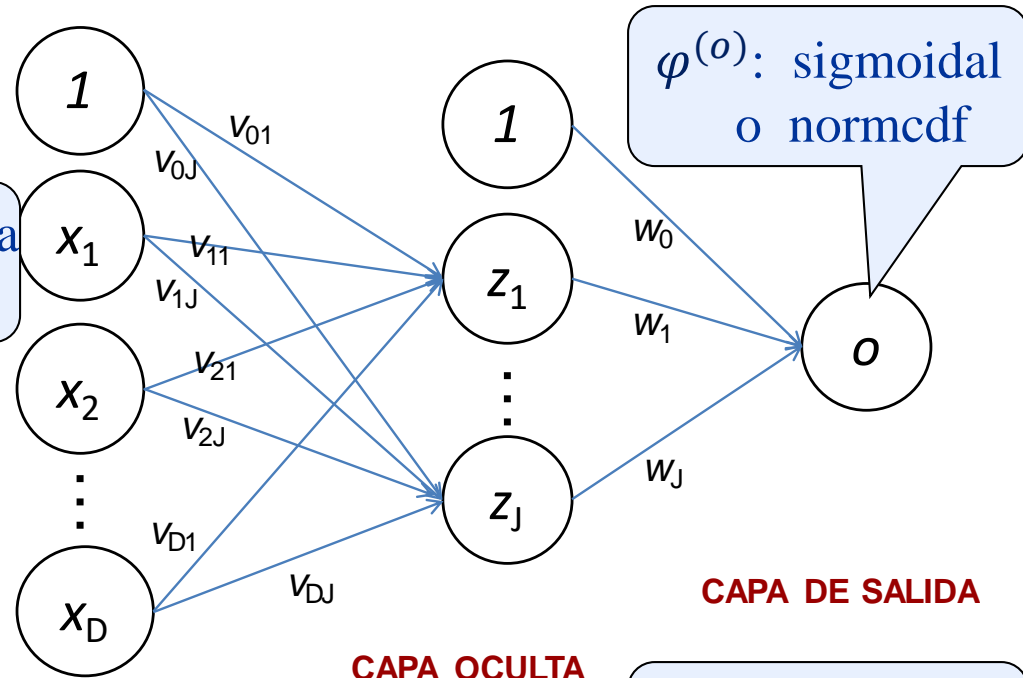
$$\mathbf{v}_j^T = (v_{0j} \ v_{1j} \ \dots \ v_{Dj});$$

$$\mathbf{z} = \begin{pmatrix} 1 \\ z_1 \\ \vdots \\ z_J \end{pmatrix};$$

Local field

$$z_j = \phi_j^{(H)}(\mathbf{v}_j^T \mathbf{x}); \quad j = 1, 2, \dots, J;$$

Función de
Activación



$$\mathbf{w}^T = (w_0 \ w_1 \ \dots \ w_J);$$

$$o(\mathbf{x}; \mathbf{w}) = \varphi^{(o)}(\mathbf{w}^T \mathbf{z}) = \varphi^{(o)}\left(\sum_{j=0}^J w_j \phi_j^{(H)}(\mathbf{v}_j^T \mathbf{x})\right)$$

Función de activación sigmoide (logit)

$$\sigma(z) = \frac{1}{1 + e^{-z}};$$

También: función de transferencia

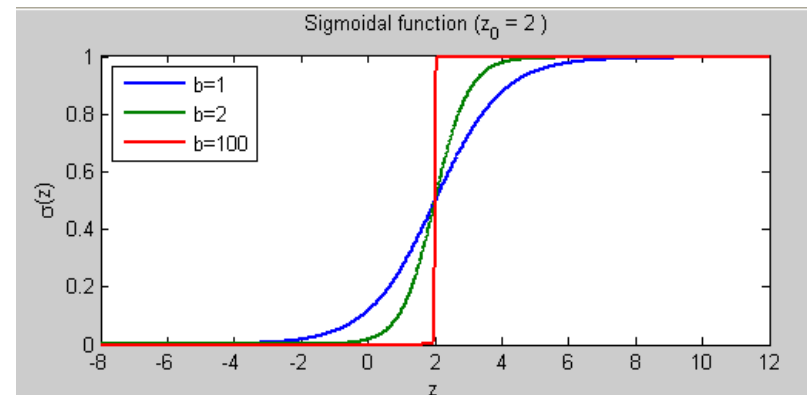
- $\sigma(-\infty) = 0; \quad \sigma(0) = \frac{1}{2}; \sigma(+\infty) = 1;$
- Monotónamente creciente : $z_2 > z_1 \Rightarrow \sigma(z_2) > \sigma(z_1)$
- Simetría : $\sigma(-z) = 1 - \sigma(z)$
- Derivada: $\sigma'(z) = \frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$

$$\sigma(z; z_0, b) = \frac{1}{1 + e^{-b(z-z_0)}}$$

$\frac{1}{b}$: escala

z_0 : centro

$$\lim_{b \rightarrow \infty} \sigma(z; z_0, b) = \begin{cases} 0 & \text{si } z < z_0 \\ \frac{1}{2} & \text{si } z = z_0 \\ 1 & \text{si } z > z_0 \end{cases}$$



Función de paso de Heaviside

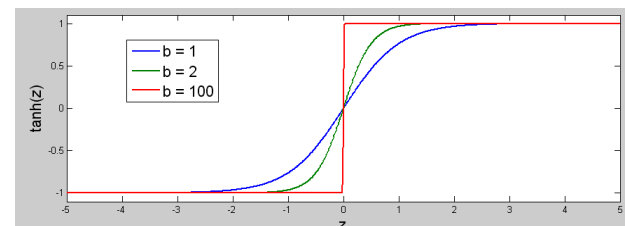
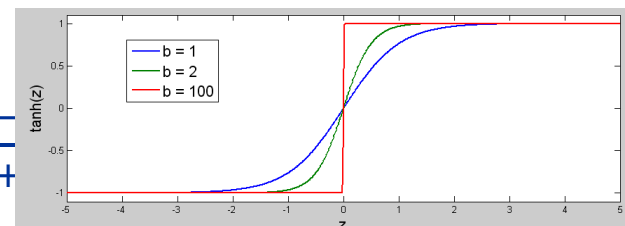
Otras funciones de activación

In regression output layer

- **Lineal:** $\varphi(z) = z$

- **Tangente hiperbolica :**

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

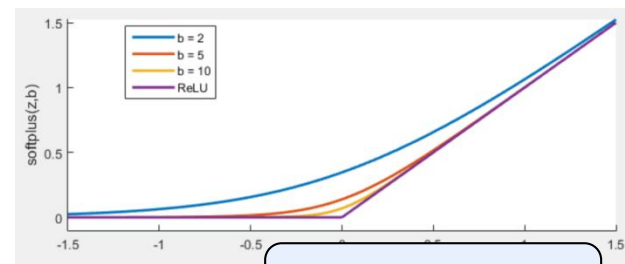


- **Probit:** $\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{y^2}{2}} dy$

normcdf

- **softplus:** $\text{softplus}(z; z_0, b) = \frac{\log(1 + e^{-b(z-z_0)})}{b}$

$$\lim_{b \rightarrow \infty} \text{softplus}(z; 0, b) = \text{ReLU}(z)$$



- **Unidad lineal rectificada :** $\text{ReLU}(z) = \max(0, z)$

Ej. $a = 0.01$

- **Unidad lineal rectificada con fugas :** $\text{Leaky ReLU}(z) = \begin{cases} az & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases}; \quad 0 < a \ll 1$

Evita el gradiente cero si $z < 0$

Perceptrón multicapa: aprendizaje en línea

ENTRADA:

Training instances:

$$\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^N$$

Parámetro de aprendizaje : $\eta > 0$

SALIDA:

$$\mathbf{V}^*, \mathbf{w}^* = \arg \min_{\mathbf{V}, \mathbf{w}} CE(\mathbf{V}, \mathbf{w})$$

Solo una capa oculta.
Activaciones sigmoideas.

1. Inicializar aleatoriamente $\mathbf{V}, \mathbf{w} \sim U[-0.5, 0.5]^{(D+1)}$

2. $n_{epoch} = 0$

3. Mientras no se cumplan los criterios de convergencia

3.1 Incrementar el contador de épocas : $n_{epoch} = n_{epoch} + 1$

3.2 Para $n = 1, \dots, N$

$$z_{nj} = \sigma(\mathbf{v}_j^T \mathbf{x}_n), \quad j = 1, 2, \dots, J; \quad \# \text{ propagación hacia adelante}$$

$$o_n = \sigma(\mathbf{w}^T \mathbf{z}_n) \quad \# \text{ salida de red}$$

$$\delta_n = o_n - t_n; \quad \# \text{ error de predicción}$$

$$\mathbf{w} = \mathbf{w} - \eta \delta_n \mathbf{z}_n \quad \# \text{ actualización de peso}$$

$$\Delta_{nj} = z_{nj}(1 - z_{nj})\delta_n, \quad j = 1, 2, \dots, J; \quad \# \text{ error de propagación hacia atrás}$$

$$\mathbf{v}_j = \mathbf{v}_j - \eta \Delta_{nj} \mathbf{x}_n \quad \# \text{ actualización de peso}$$

Propiedad de aproximación universal

THEOREM: “A feed-forward network with a **single hidden layer** containing a **sufficiently large number of hidden neurons** can uniformly **approximate any continuous function** on compact subsets of \mathbb{R}^D , under mild conditions on the activation function”

■ GOOD NEWS:

A neural network can be used to make optimal predictions!

■ NOT SO GOOD NEWS:

How does one find such a network?

- How many neurons do we need?
- How easy is it to learn the network parameters (weights)?

Use CV

Having more than one layer can help (deep networks)

K. Hornik, M. Stinchcombe and H. White, “Multi-layer Feedforward Networks are Universal Approximators,” Neural Networks, Vol. 2, pp. 359-366 (1989).

Propiedad de aproximación universal

TEOREMA: “Una red de retroalimentación con una sola capa oculta que contiene un **número suficientemente grande de neuronas ocultas** puede **aproximar** de manera uniforme **cualquier función continua** en subconjuntos compactos de \mathbb{R}^D , bajo condiciones suaves en la función de activación”

■ BUENAS NOTICIAS :

¡Se puede usar una red neuronal para hacer predicciones óptimas!

■ NO TAN BUENAS NOTICIAS :

¿Cómo se encuentra una red así?

- ¿Cuántas neuronas necesitamos?
- ¿Qué tan fácil es aprender los parámetros de red (pesos)?

Usar VC

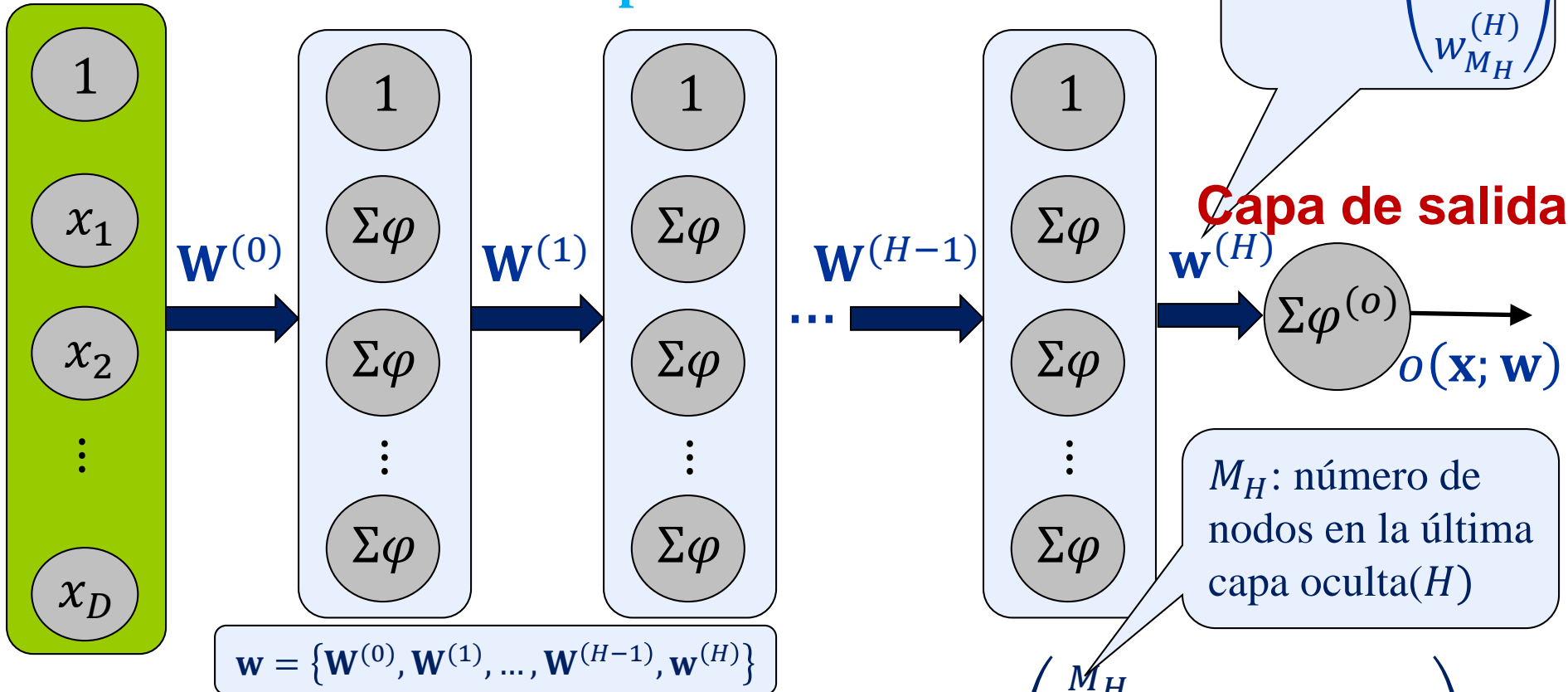
Tener más de una capa puede ayudar (redes profundas)

K. Hornik, M. Stinchcombe and H. White, “Multi-layer Feedforward Networks are Universal Approximators,” Neural Networks, Vol. 2, pp. 359-366 (1989).

Redes neuronales profundas: clasificación

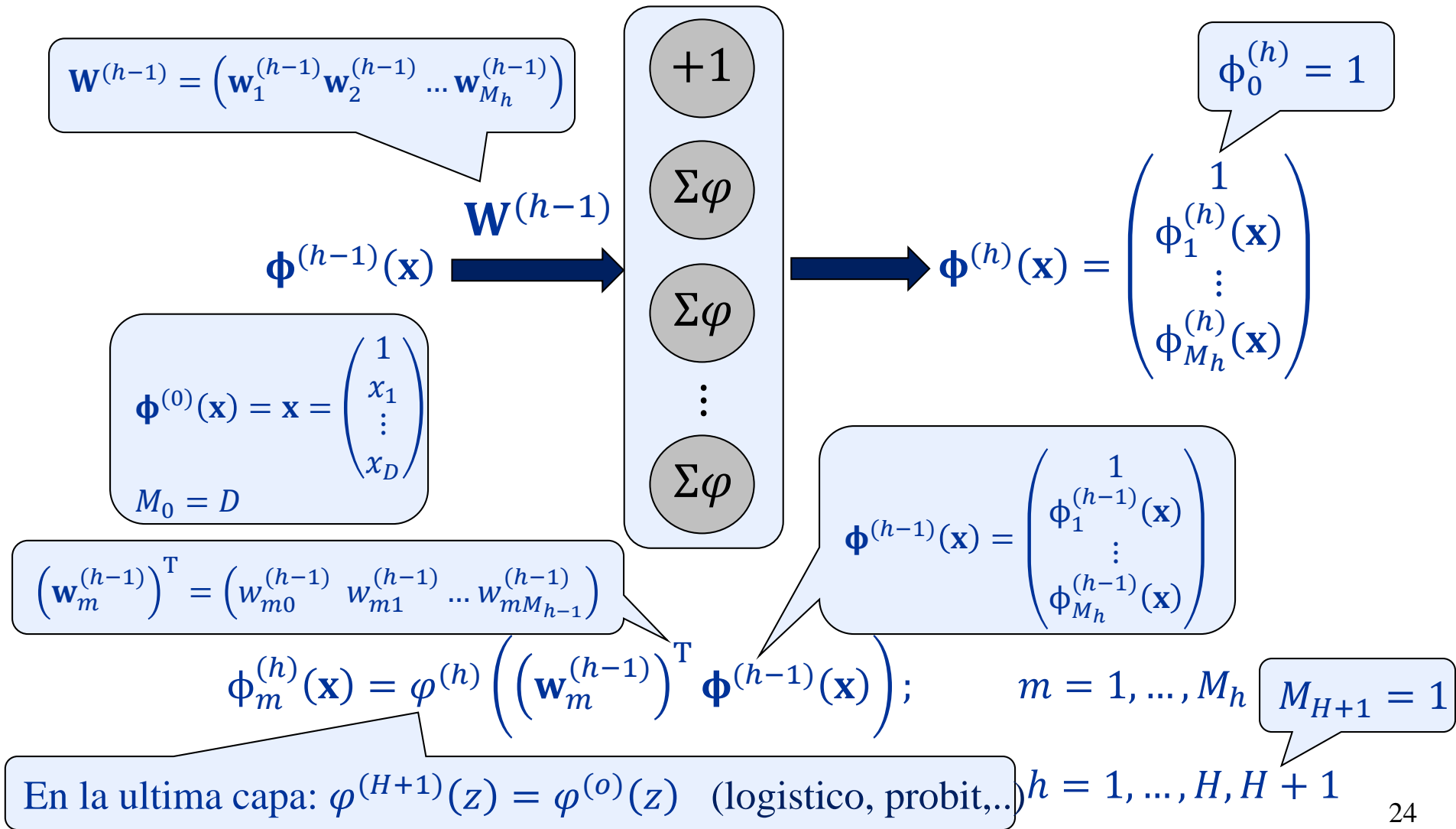
Capa de entrada

H capas ocultas



$$o(\mathbf{x}; \mathbf{w}) = \varphi^{(o)} \left((\mathbf{w}^{(H)})^T \boldsymbol{\Phi}^{(H)}(\mathbf{x}) \right) = \varphi^{(o)} \left(\sum_{m=0}^{M_H} w_m^{(H)} \phi_m^{(H)}(\mathbf{x}) \right)$$

Salida de capa oculta h



Aprendiendo los pesos de la red

Los pesos de la red se determinan minimizando una función de coste:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left(\left[- \sum_{n \in C_0} \log(1 - o(\mathbf{x}_n; \mathbf{w})) - \sum_{n \in C_1} \log(o(\mathbf{x}_n; \mathbf{w})) \right] + \lambda_1 \|\mathbf{w}\|_1 + \lambda_2 \|\mathbf{w}\|_2^2 \right)$$

Uso VC para seleccionar

$\lambda_1 > 0$ L_1 penalización

$\lambda_2 > 0$ L_2 penalización

Error de entropía cruzada

- **Arquitectura:** número de capas ocultas / neuronas en capa oculta
- **Método de optimización:** cuasi-Newton, descenso de gradiente, descenso de gradiente estocástico, uso de impulso, ... Gradiientes necesarios!
- **Hiperparámetros:** magnitud de las penalizaciones, del plazo de impulso, # máximo de iteraciones, ... Impulso/momentum es acelerar la velocidad de aprendizaje

Clasificación de clases múltiples (clases K)

Codificación 1-de K de los datos etiquetados

$$\mathcal{D} = \{(\mathbf{x}_n, c_n)\}_{n=1}^N \Rightarrow \mathcal{D} = \{(\mathbf{x}_n, \mathbf{t}_n)\}_{n=1}^N; \quad \mathbf{t}_n^T = (t_{n1} \ t_{n2} \ \dots \ t_{nK})$$

$$t_{nk} = \mathbb{I}[c_n = C_k] = \begin{cases} 0 & \text{if } y_n \neq C_k; \\ 1 & \text{if } y_n = C_k; \end{cases} \quad k = 1, 2, \dots, K$$

$$\mathcal{L}: \mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N \rightarrow o(\cdot; \mathbf{w})$$

$$\mathbf{z}(\cdot; \mathbf{w}): \mathbf{x} \in \mathcal{X} \rightarrow \mathbf{z}(\mathbf{x}; \mathbf{w}) \in \mathbb{R}^K$$

$$\mathbf{o}(\mathbf{x}; \mathbf{w}) = \boldsymbol{\varphi}^{(o)}(\mathbf{z}(\mathbf{x}; \mathbf{w})) \in \Delta^K$$

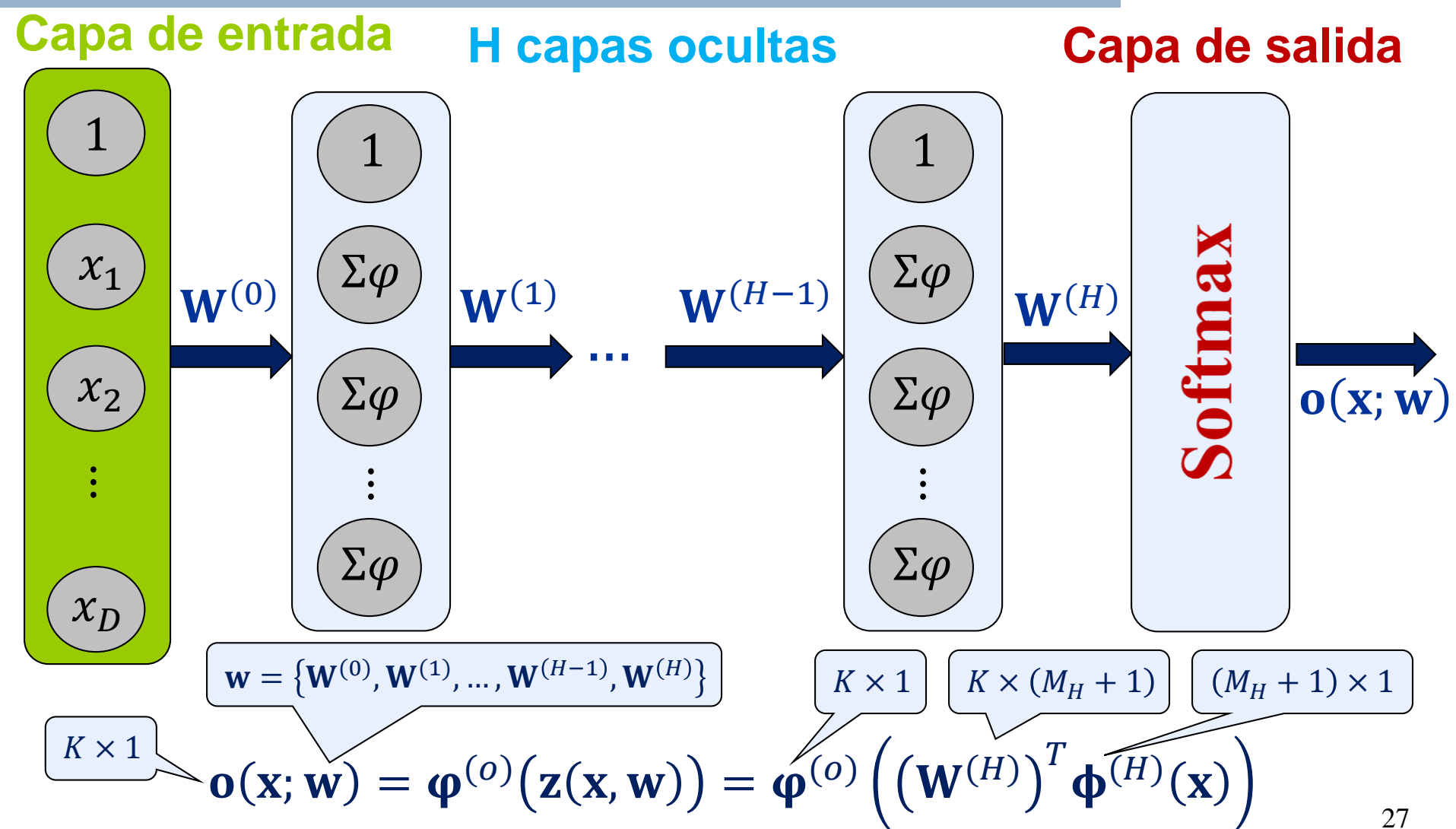
model
Discriminativo
 $\hat{p}(C_k | \mathbf{x}, \mathbf{w}) = \varphi_K^{(o)}(\mathbf{z}(\mathbf{x}; \mathbf{w}))$

Probabilidad simplex
en K dimensiones

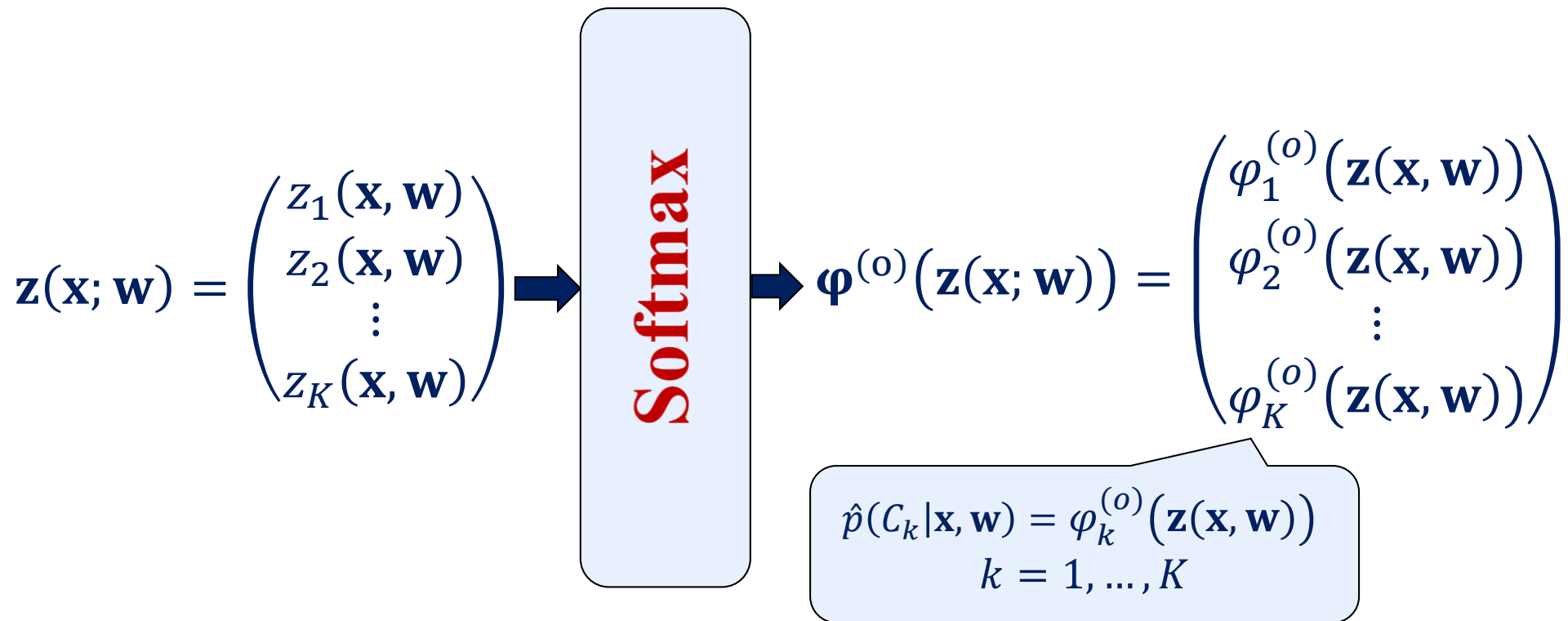
$$\left(\boldsymbol{\varphi}^{(o)}(\mathbf{z})\right)^T = \left(\varphi_1^{(o)}(\mathbf{z}) \ \dots \ \varphi_K^{(o)}(\mathbf{z})\right)$$

$$\begin{aligned} \varphi_k^{(o)}(\mathbf{z}) &\geq 0 & k = 1, \dots, K \\ \varphi_1^{(o)}(\mathbf{z}) + \dots + \varphi_K^{(o)}(\mathbf{z}) &= 1 \end{aligned}$$

MLP para clasificación multiclase



Capa Softmax



probabilidades

$$\varphi_k^{(o)}(\mathbf{z}) = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}} \geq 0 \quad k = 1, \dots, K$$

$$\varphi_1^{(o)}(\mathbf{z}) + \dots + \varphi_K^{(o)}(\mathbf{z}) = 1$$

Aprendizaje por máxima probabilidad

Se asumen
Muestras dii
(distribuidas
Idéntica e
Independiente-
mente)

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \mathcal{L}(\mathbf{w})$$

- Probabilidad del modelo, dado $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{t}_n)\}_{n=1}^N$

$$\mathcal{L}(\mathbf{w}) = \hat{P}(\{\mathbf{t}_n\}_{n=1}^N | \{\mathbf{x}_n\}_{n=1}^N, \mathbf{w}) = \prod_{n=1}^N \hat{p}(t_n | \mathbf{x}_n, \mathbf{w}) =$$

Factoriza porque
se asumen muestras
independientes

$$= \prod_{n=1}^N \prod_{k=1}^K (\hat{p}(C_k | \mathbf{x}_n, \mathbf{w}))^{t_{nk}}$$

Igualamos factores porque
se asumen muestras
idénticamente dist.

- Función de verosimilitud-log

$$\log \mathcal{L}(\mathbf{w}) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log \hat{p}(C_k | \mathbf{x}_n, \mathbf{w})$$

$$t_{nk} = \mathbb{I}[c_n = C_k]$$

Minimización del error de entropía cruzada

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} CE(\mathbf{w})$$

- Error de entropía cruzada

$$CE(\mathbf{w}) = -\log \mathcal{L}(\mathbf{w}) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \log \hat{p}(C_k | \mathbf{x}_n, \mathbf{w})$$

- Modelo discriminatorio: $\hat{p}(C_1 | \mathbf{x}, \mathbf{w}) = \varphi_k^{(o)}(z(\mathbf{x}_n, \mathbf{w}))$

Salida de la
k neurona en
la última capa
de la MLP.

$$\widehat{CE}(\mathbf{w}) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \log \varphi_k^{(o)}(z(\mathbf{x}_n, \mathbf{w}))$$

Resumen de redes neuronales

■ Ventajas

- Excelentes predictores.
- Adaptativo (aprendizaje en línea)

■ Desventajas

- Entrenamiento costoso.
- Encontrar la arquitectura adecuada puede ser difícil
 - Número de capas / nodos ocultos en cada capa oculta
 - Función de activación
- Determinando los hiperparámetros para la optimización
 - Tipo de optimización
 - En SGD: tasa de aprendizaje, tamaño de mini lotes, término de impulso, fuerza de los términos de regularización, ...
- Difícil interpretación.

Estado del arte: redes neuronales profundas

Resumen de redes neuronales

■ Advantages

- Excellent predictors.
- Adaptive (online learning)

■ Disadvantages

- Costly training.
- Finding appropriate architecture can be difficult
 - Number of hidden layers / nodes in each hidden layer
 - Activation function
- Determining the hyperparameters for the optimization
 - Type of optimization
 - In SGD: learning rate, size of mini-batches, momentum term, strength of regularization terms, ...
- Difficult interpretation.

State-of-the-art: Deep neural networks

Referencias

- Multilayer perceptron with TensorFlow
<https://playground.tensorflow.org/>
- Scikit learn documentation
 - Neural networks
http://scikit-learn.org/stable/modules/neural_networks_supervised.html
- Simon O. Haykin “Neural Networks and Learning Machines, 3rd edition”, Pearson (2009)
- Ian Goodfellow and Yoshua Bengio and Aaron Courville “Deep Learning”, MIT Press (2016)