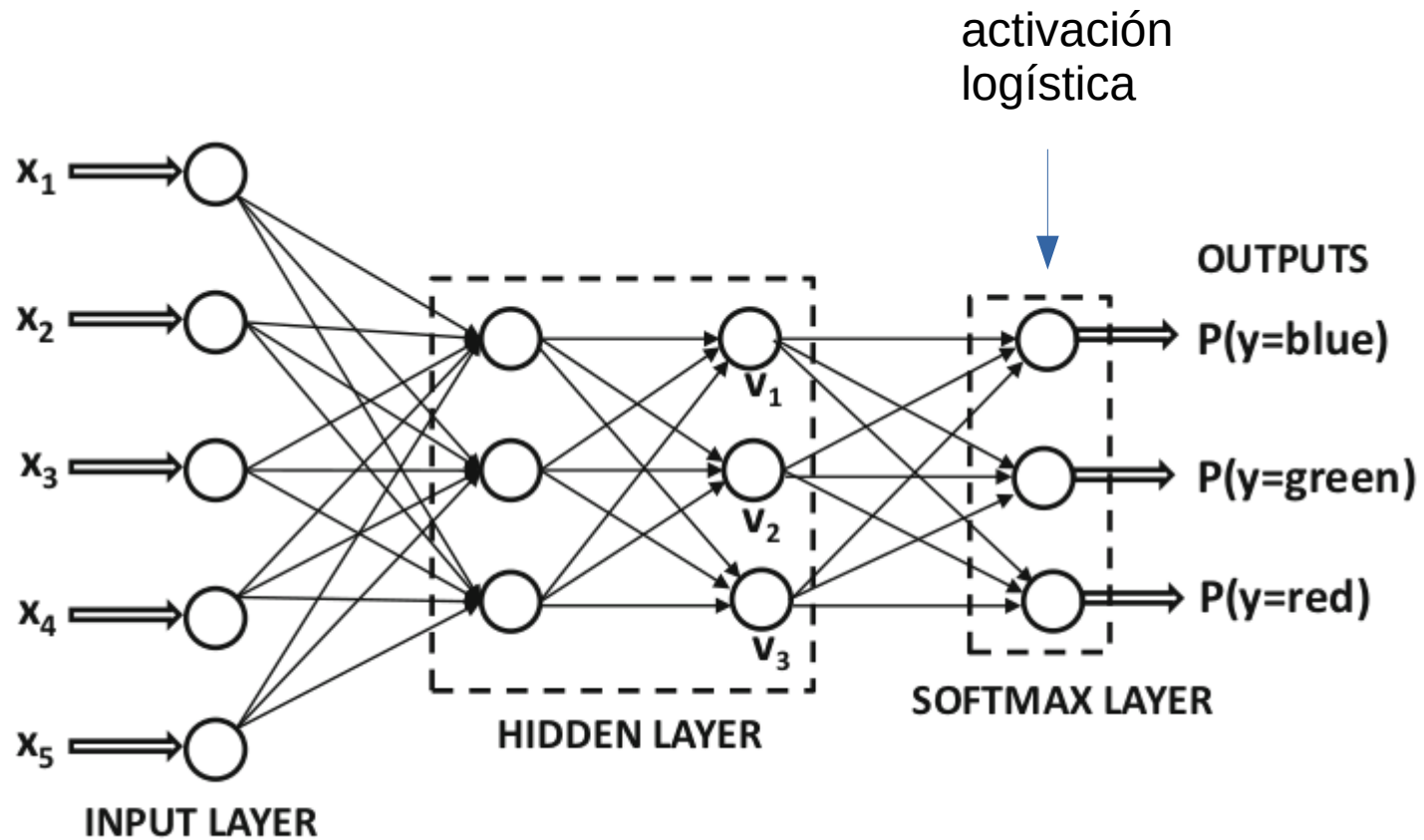


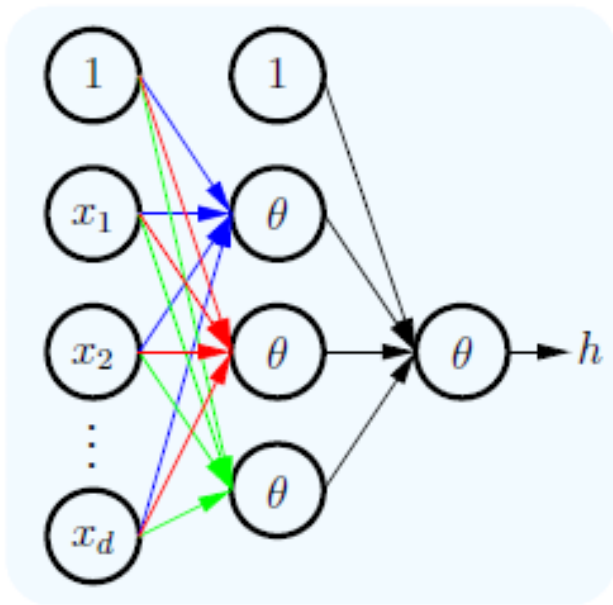


# Minería de Datos

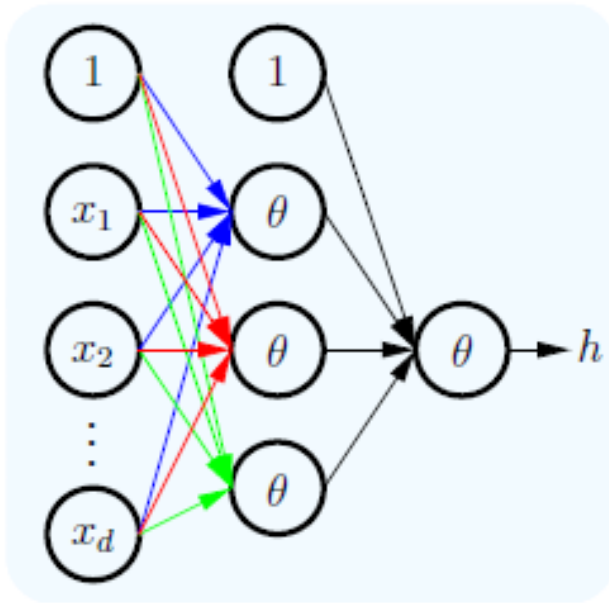
## Extensión de MLP a $K$ clases



## Funciones aproximadas por una red feed-forward

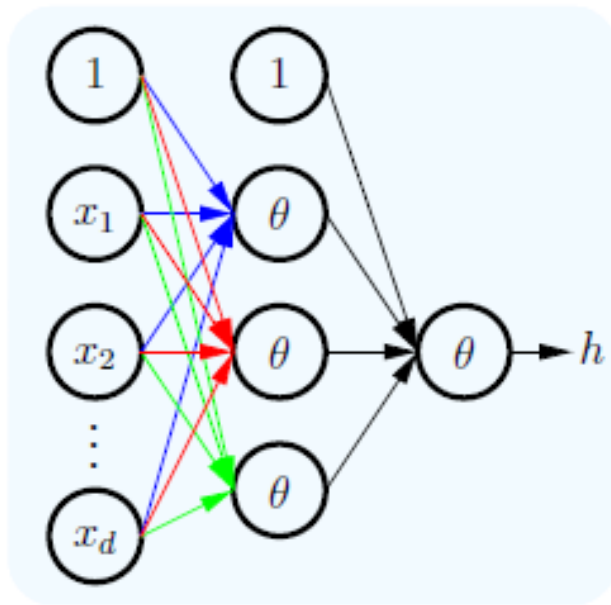


## Funciones aproximadas por una red feed-forward



$$h(\mathbf{x}) = \theta \left( w_{01}^{(2)} + \sum_{j=1}^m w_{j1}^{(2)} \theta \left( \sum_{i=0}^d w_{ij}^{(1)} x_i \right) \right)$$

## Funciones aproximadas por una red feed-forward



$$h(\mathbf{x}) = \theta \left( w_{01}^{(2)} + \sum_{j=1}^m w_{j1}^{(2)} \theta \left( \sum_{i=0}^d w_{ij}^{(1)} x_i \right) \right)$$

Se usa también la siguiente notación (simplificada):

$$h(\mathbf{x}) = \theta \left( w_0 + \sum_{j=1}^m w_j \theta \left( \mathbf{v}_j^T \mathbf{x} \right) \right)$$

$\nearrow W^{(1)}$

Y para dimensionalidad:  $d^{(1)} = m$

## Capacidad de una red feed-forward

Para una cantidad de neuronas suficientemente grande, podemos lograr:

$$E_{\text{in}} \sim 0$$

Sin embargo, no podemos asegurar que  $E_{\text{in}} \approx E_{\text{out}}$ .

## Capacidad de una red feed-forward

Para una cantidad de neuronas suficientemente grande, podemos lograr:

$$E_{\text{in}} \sim 0$$

Sin embargo, no podemos asegurar que  $E_{\text{in}} \approx E_{\text{out}}$ .

Se puede mostrar que para esta red:

$$h(\mathbf{x}) = \theta \left( w_{01}^{(2)} + \sum_{j=1}^m w_{j1}^{(2)} \theta \left( \sum_{i=0}^d w_{ij}^{(1)} x_i \right) \right)$$

la dimensión VC está acotada por:

$$d_{\text{VC}} \leq (\text{const}) \cdot md \log(md).$$

si la activación es la función signo en ambas capas.

## Capacidad de una red feed-forward

Si la red usa estas funciones de activación:

$$h(\mathbf{x}) = \underset{\substack{\uparrow \\ \text{sign}(x)}}{\theta} \left( w_{01}^{(2)} + \sum_{j=1}^m w_{j1}^{(2)} \underset{\substack{\uparrow \\ \tanh(\cdot)}}{\theta} \left( \sum_{i=0}^d w_{ij}^{(1)} x_i \right) \right)$$

la dimensión VC es:

$$d_{\text{VC}} = O(md(m + d)).$$

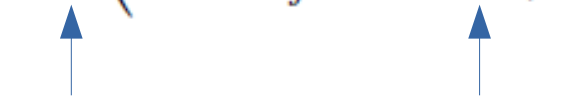
¿Cuál de las dos redes tiene mayor expresividad?



## Capacidad de una red feed-forward

Si la red usa estas funciones de activación:

$$h(\mathbf{x}) = \theta \left( w_{01}^{(2)} + \sum_{j=1}^m w_{j1}^{(2)} \theta \left( \sum_{i=0}^d w_{ij}^{(1)} x_i \right) \right)$$

  
 $\text{sign}(x)$                        $\tanh(\cdot)$

la dimensión VC es:

$$d_{\text{VC}} = O(md(m + d)).$$

¿Cuál de las dos redes tiene mayor expresividad?

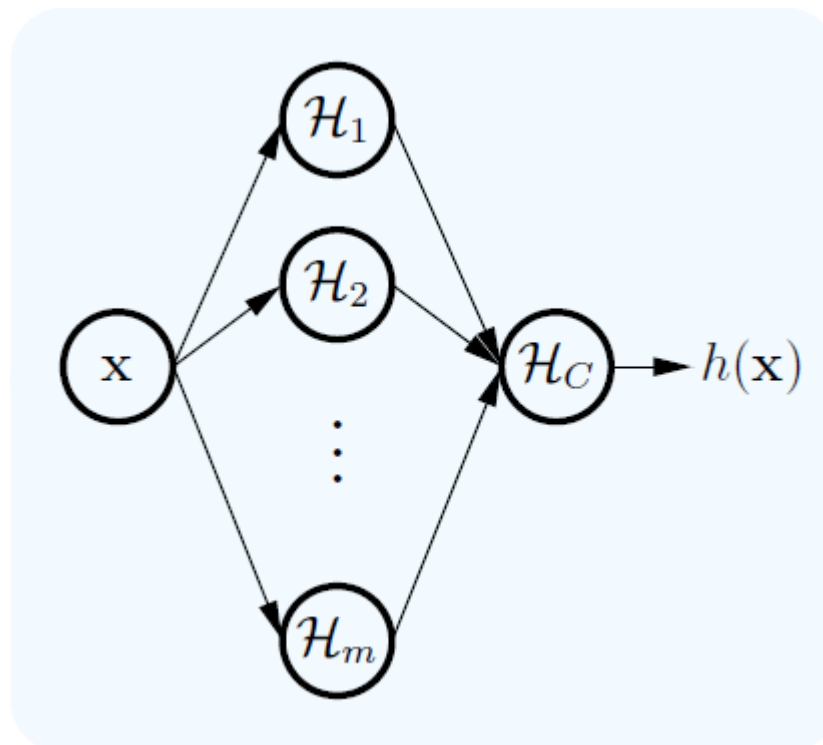
$$md \log(md) \leq md(m + d)$$

Usar  $\tanh(\cdot)$  aumenta la expresividad de la red.

## Capacidad de una red feed-forward

¿De dónde viene  $d_{VC} \leq (\text{const}) \cdot md \log(md)$ .?

Consideremos el siguiente conjunto de hipótesis:



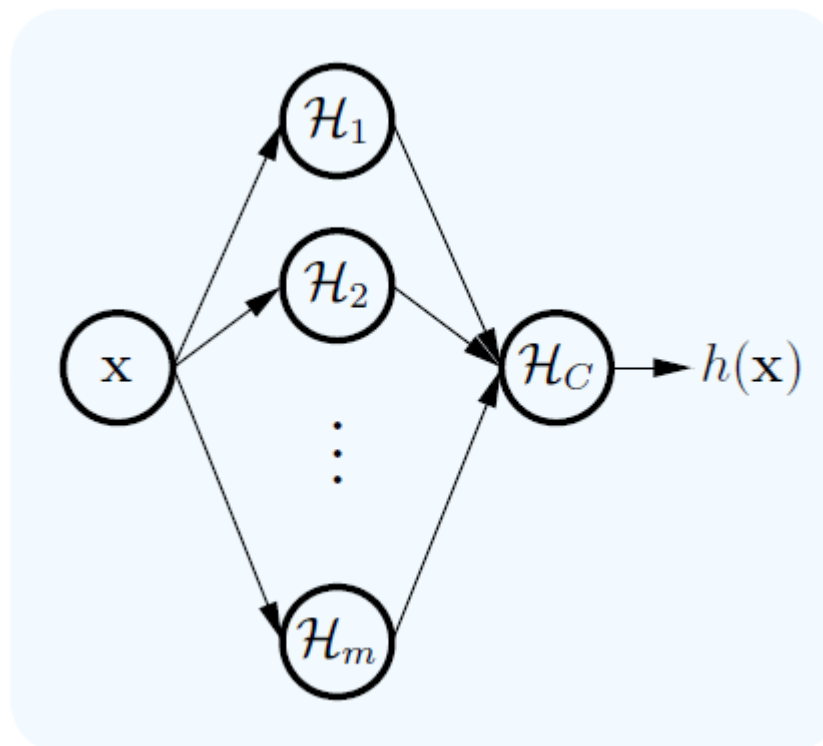
$$h(\mathbf{x}) = h_C(h_1(\mathbf{x}), \dots, h_m(\mathbf{x}))$$

↑  
perceptrones

## Capacidad de una red feed-forward

¿De dónde viene  $d_{VC} \leq (\text{const}) \cdot md \log(md)$ .?

Consideremos el siguiente conjunto de hipótesis:



$$h(\mathbf{x}) = h_C(h_1(\mathbf{x}), \dots, h_m(\mathbf{x}))$$

↑  
perceptrones

Supongamos que la dimensión VC de  $\mathcal{H}_i$  es  $d_i$  y la de  $\mathcal{H}_C$  es  $d_c$ .

## Capacidad de una red feed-forward

Fijemos  $\mathbf{x}_1, \dots, \mathbf{x}_N$  y  $h_1, \dots, h_m$ . Las hipótesis son ahora funciones base que definen una transformación hacia  $\mathbb{R}^m$ ,

$$\mathbf{x}_1 \rightarrow \mathbf{z}_1 = \begin{bmatrix} h_1(\mathbf{x}_1) \\ \vdots \\ h_m(\mathbf{x}_1) \end{bmatrix} \quad \dots \quad \rightarrow \mathbf{z}_N = \begin{bmatrix} h_1(\mathbf{x}_N) \\ \vdots \\ h_m(\mathbf{x}_N) \end{bmatrix} .$$

## Capacidad de una red feed-forward

Fijemos  $\mathbf{x}_1, \dots, \mathbf{x}_N$  y  $h_1, \dots, h_m$ . Las hipótesis son ahora funciones base que definen una transformación hacia  $\mathbb{R}^m$ ,

$$\mathbf{x}_1 \rightarrow \mathbf{z}_1 = \begin{bmatrix} h_1(\mathbf{x}_1) \\ \vdots \\ h_m(\mathbf{x}_1) \end{bmatrix} \quad \dots \quad \rightarrow \mathbf{z}_N = \begin{bmatrix} h_1(\mathbf{x}_N) \\ \vdots \\ h_m(\mathbf{x}_N) \end{bmatrix}.$$

Los vectores  $\mathbf{z}$  son vectores binarios en  $\mathbb{R}^m$ .

Dado que la red tiene flexibilidad para encontrar  $h_1, \dots, h_m$ , podemos acotar por arriba el número posible de diferentes  $\mathbf{z}_1, \dots, \mathbf{z}_N$ .

## Capacidad de una red feed-forward

Fijemos  $\mathbf{x}_1, \dots, \mathbf{x}_N$  y  $h_1, \dots, h_m$ . Las hipótesis son ahora funciones base que definen una transformación hacia  $\mathbb{R}^m$ ,

$$\mathbf{x}_1 \rightarrow \mathbf{z}_1 = \begin{bmatrix} h_1(\mathbf{x}_1) \\ \vdots \\ h_m(\mathbf{x}_1) \end{bmatrix} \quad \dots \quad \rightarrow \mathbf{z}_N = \begin{bmatrix} h_1(\mathbf{x}_N) \\ \vdots \\ h_m(\mathbf{x}_N) \end{bmatrix}.$$

Los vectores  $\mathbf{z}$  son vectores binarios en  $\mathbb{R}^m$ .

Dado que la red tiene flexibilidad para encontrar  $h_1, \dots, h_m$ , podemos acotar por arriba el número posible de diferentes  $\mathbf{z}_1, \dots, \mathbf{z}_N$ .

Las primeras componentes de los vectores  $\mathbf{z}$  están dadas por:

$$h_1(\mathbf{x}_1), \dots, h_1(\mathbf{x}_N)$$

## Capacidad de una red feed-forward

Fijemos  $\mathbf{x}_1, \dots, \mathbf{x}_N$  y  $h_1, \dots, h_m$ . Las hipótesis son ahora funciones base que definen una transformación hacia  $\mathbb{R}^m$ ,

$$\mathbf{x}_1 \rightarrow \mathbf{z}_1 = \begin{bmatrix} h_1(\mathbf{x}_1) \\ \vdots \\ h_m(\mathbf{x}_1) \end{bmatrix} \quad \dots \quad \rightarrow \mathbf{z}_N = \begin{bmatrix} h_1(\mathbf{x}_N) \\ \vdots \\ h_m(\mathbf{x}_N) \end{bmatrix}.$$

Los vectores  $\mathbf{z}$  son vectores binarios en  $\mathbb{R}^m$ .

Dado que la red tiene flexibilidad para encontrar  $h_1, \dots, h_m$ , podemos acotar por arriba el número posible de diferentes  $\mathbf{z}_1, \dots, \mathbf{z}_N$ .

Las primeras componentes de los vectores  $\mathbf{z}$  están dadas por:

$$h_1(\mathbf{x}_1), \dots, h_1(\mathbf{x}_N)$$

que es una dicotomía sobre  $\mathbf{x}_1, \dots, \mathbf{x}_N$  implementada por  $h_1$ . Dado que la dimensión VC de  $\mathcal{H}_1$  es  $d_1$ , existen a lo más  $N^{d_1}$  dicotomías. Es decir, existen a lo más  $N^{d_1}$  formas de elegir las primeras componentes de  $\mathbf{z}$ .

## Capacidad de una red feed-forward

Fijemos  $\mathbf{x}_1, \dots, \mathbf{x}_N$  y  $h_1, \dots, h_m$ . Las hipótesis son ahora funciones base que definen una transformación hacia  $\mathbb{R}^m$ ,

$$\mathbf{x}_1 \rightarrow \mathbf{z}_1 = \begin{bmatrix} h_1(\mathbf{x}_1) \\ \vdots \\ h_m(\mathbf{x}_1) \end{bmatrix} \quad \dots \quad \rightarrow \mathbf{z}_N = \begin{bmatrix} h_1(\mathbf{x}_N) \\ \vdots \\ h_m(\mathbf{x}_N) \end{bmatrix}.$$

Los vectores  $\mathbf{z}$  son vectores binarios en  $\mathbb{R}^m$ .

Dado que la red tiene flexibilidad para encontrar  $h_1, \dots, h_m$ , podemos acotar por arriba el número posible de diferentes  $\mathbf{z}_1, \dots, \mathbf{z}_N$ .

Las primeras componentes de los vectores  $\mathbf{z}$  están dadas por:

$$h_1(\mathbf{x}_1), \dots, h_1(\mathbf{x}_N)$$

que es una dicotomía sobre  $\mathbf{x}_1, \dots, \mathbf{x}_N$  implementada por  $h_1$ . Dado que la dimensión VC de  $\mathcal{H}_1$  es  $d_1$ , existen a lo más  $N^{d_1}$  dicotomías. Es decir, existen a lo más  $N^{d_1}$  formas de elegir las primeras componentes de  $\mathbf{z}$ .

—► Recordar que:



## Capacidad de una red feed-forward

Luego, el número total de asignaciones para  $\mathbf{z}_1, \dots, \mathbf{z}_N$  es:

$$\prod_{i=1}^m N^{d_i} = N^{\sum_{i=1}^m d_i}$$

## Capacidad de una red feed-forward

Luego, el número total de asignaciones para  $\mathbf{z}_1, \dots, \mathbf{z}_N$  es:

$$\prod_{i=1}^m N^{d_i} = N^{\sum_{i=1}^m d_i}$$

Cada una de estas asignaciones puede ser dicotomizada en a lo más  $N^{d_c}$  formas. Cada una de estas dicotomías para una asignación particular entrega una dicotomía para los datos  $\mathbf{x}_1, \dots, \mathbf{x}_N$ . Luego, el número máximo de dicotomías implementables sobre los datos es acotada superiormente por el producto:

$$m(N) \leq N^{d_c} \cdot N^{\sum_{i=1}^m d_i} = N^{d_c + \sum_{i=1}^m d_i}.$$

## Capacidad de una red feed-forward

Luego, el número total de asignaciones para  $\mathbf{z}_1, \dots, \mathbf{z}_N$  es:

$$\prod_{i=1}^m N^{d_i} = N^{\sum_{i=1}^m d_i}$$

Cada una de estas asignaciones puede ser dicotomizada en a lo más  $N^{d_c}$  formas. Cada una de estas dicotomías para una asignación particular entrega una dicotomía para los datos  $\mathbf{x}_1, \dots, \mathbf{x}_N$ . Luego, el número máximo de dicotomías implementables sobre los datos es acotada superiormente por el producto:

$$m(N) \leq N^{d_c} \cdot N^{\sum_{i=1}^m d_i} = N^{d_c + \sum_{i=1}^m d_i}.$$

Para un MLP de dos capas:

$$d_i = d + 1$$

$$d_c = m + 1$$

$$\rightarrow D = d_c + \sum_{i=1}^m d_i = m(d + 2) + 1 = O(md)$$

## Capacidad de una red feed-forward

Luego, el número total de asignaciones para  $\mathbf{z}_1, \dots, \mathbf{z}_N$  es:

$$\prod_{i=1}^m N^{d_i} = N^{\sum_{i=1}^m d_i}$$

Cada una de estas asignaciones puede ser dicotomizada en a lo más  $N^{d_c}$  formas. Cada una de estas dicotomías para una asignación particular entrega una dicotomía para los datos  $\mathbf{x}_1, \dots, \mathbf{x}_N$ . Luego, el número máximo de dicotomías implementables sobre los datos es acotada superiormente por el producto:

$$m(N) \leq N^{d_c} \cdot N^{\sum_{i=1}^m d_i} = N^{d_c + \sum_{i=1}^m d_i}.$$

Para un MLP de dos capas:

$$d_i = d + 1$$

$$d_c = m + 1$$

Recordar que:

$$m_{\mathcal{H}}(N) \leq N^{d_{vc}} + 1.$$

→ Al menos separa  $\Omega(md)$  puntos.

$$\rightarrow D = d_c + \sum_{i=1}^m d_i = m(d + 2) + 1 = O(md)$$

## Capacidad de una red feed-forward

El análisis anterior puede extenderse para encontrar una cota superior ajustada:

$$d_{VC} \leq 2D \log_2 D \quad (\text{propuesto})$$

Para la MLP de dos capas, entonces:

$$d_{VC} = O(md \log(md))$$

## Capacidad de una red feed-forward

El análisis anterior puede extenderse para encontrar una cota superior ajustada:

$$d_{VC} \leq 2D \log_2 D \quad (\text{propuesto})$$

Para la MLP de dos capas, entonces:

$$d_{VC} = O(md \log(md))$$

Las neuronas de la capa oculta no deben ser demasiadas pero deben ser suficientes para ajustar el modelo a los datos.

## Capacidad de una red feed-forward

El análisis anterior puede extenderse para encontrar una cota superior ajustada:

$$d_{VC} \leq 2D \log_2 D \quad (\text{propuesto})$$

Para la MLP de dos capas, entonces:

$$d_{VC} = O(md \log(md))$$

Las neuronas de la capa oculta no deben ser demasiadas pero deben ser suficientes para ajustar el modelo a los datos.

En general, es una buena decisión que el número de neuronas crezca sublinealmente con los datos. Por ejemplo:

$$m \approx \frac{1}{d} \sqrt{N}. \rightarrow d_{VC} = O(\sqrt{N} \log \sqrt{N})$$

## Capacidad de una red feed-forward

Recordar:

$$m_{\mathcal{H}}(N) \leq N^{d_{\text{vc}}} + 1 \sim N^{d_{\text{vc}}}$$

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + O\left(\sqrt{\frac{d_{\text{vc}} \log N}{N}}\right)$$

$$\text{Si } m \approx \frac{1}{d} \sqrt{N}. \longrightarrow d_{\text{vc}} = O(\sqrt{N} \log \sqrt{N})$$

De esta manera si  $N \rightarrow \infty$ ,  $E_{\text{out}} \rightarrow E_{\text{in}}$ .

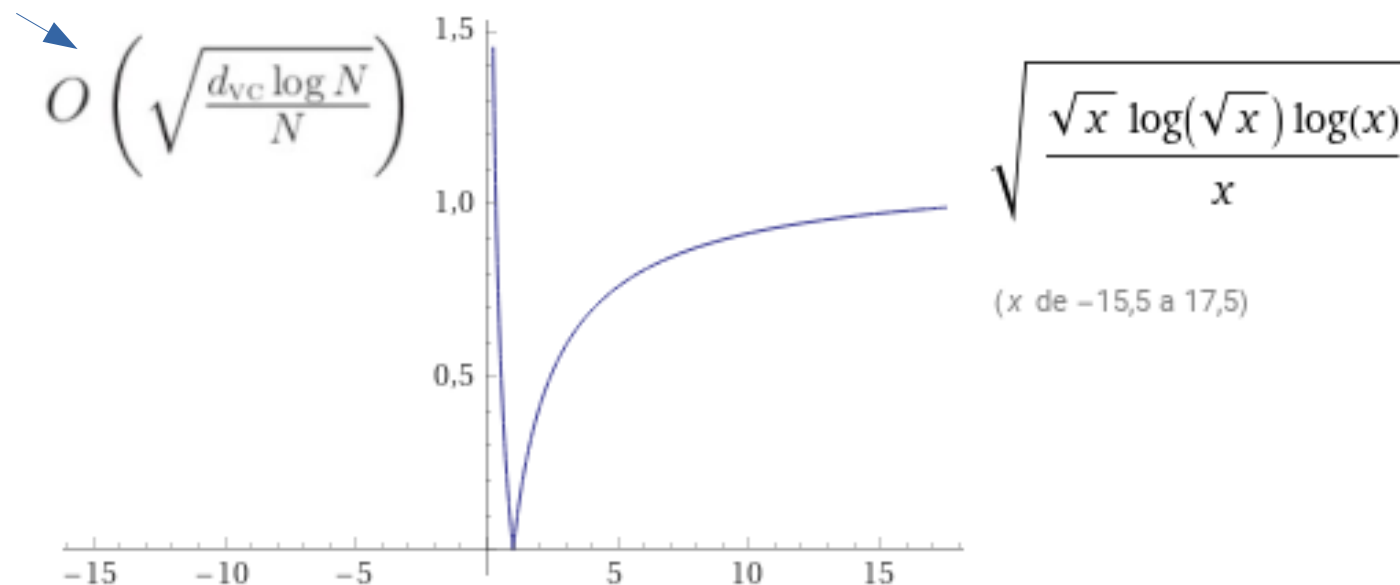


## Capacidad de una red feed-forward

$$\text{Si } m \approx \frac{1}{d} \sqrt{N}. \quad \longrightarrow \quad d_{VC} = O(\sqrt{N} \log \sqrt{N})$$

De esta manera si  $N \rightarrow \infty$ ,  $E_{\text{out}} \rightarrow E_{\text{in}}$ .

Complejidad  
paramétrica



## Aspectos prácticos en entrenamiento

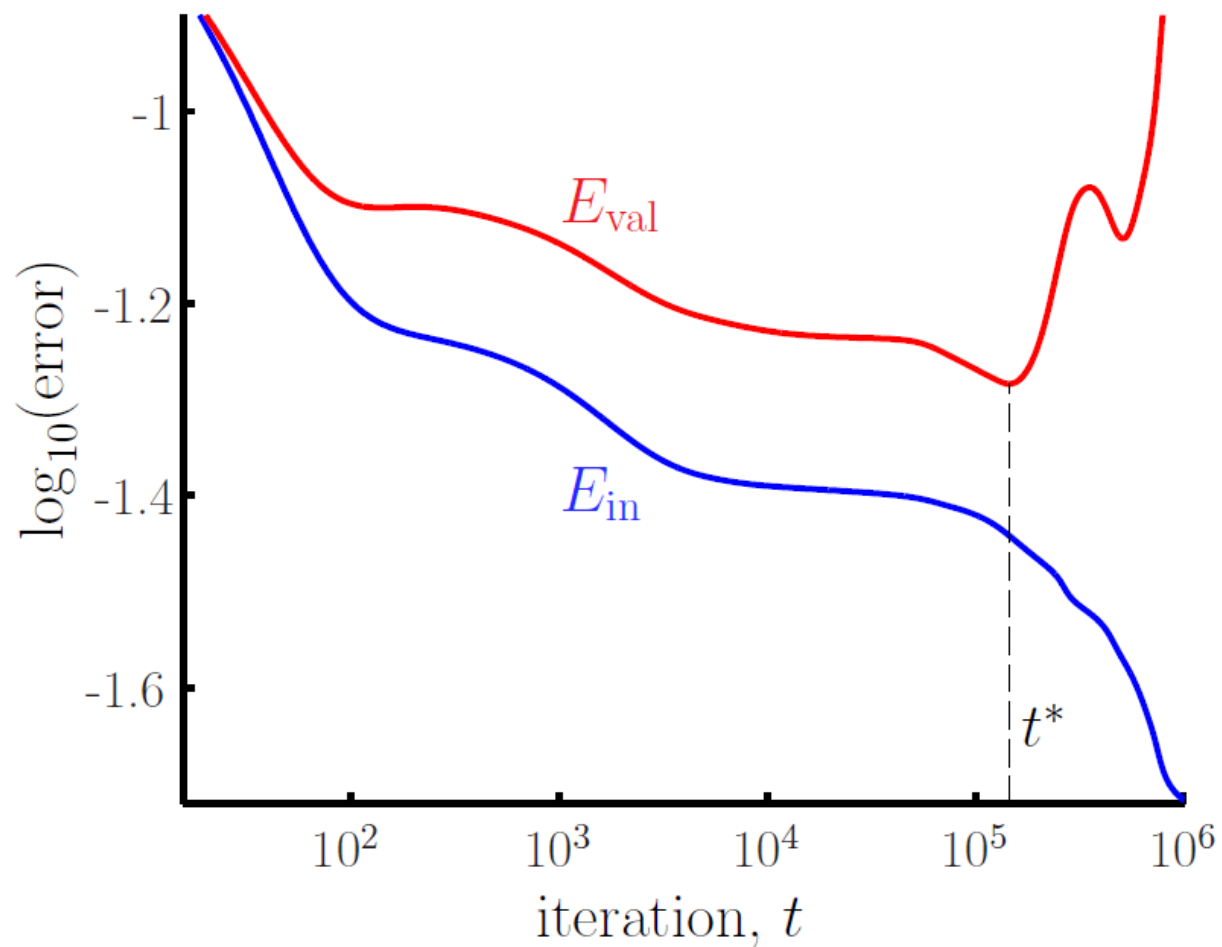
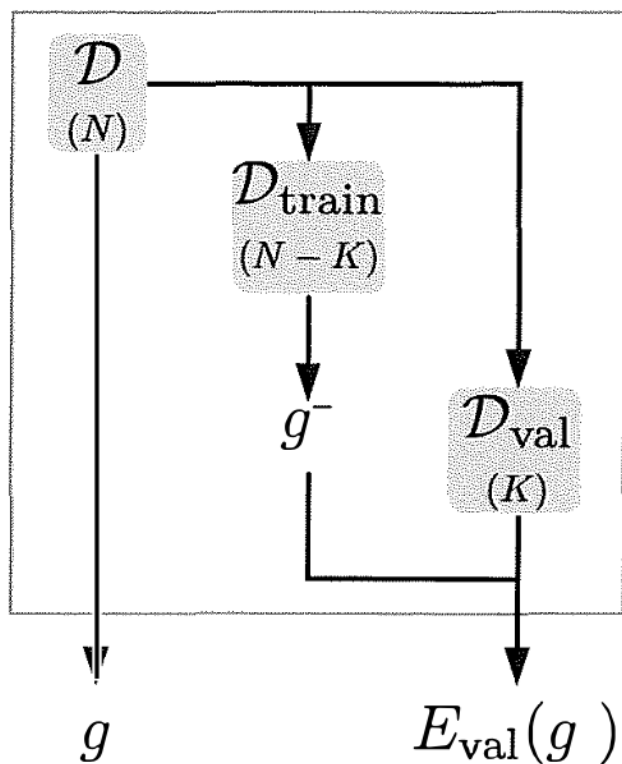
Regularización: penalizar el uso innecesario de parámetros del modelo, evitando el sobreajuste.

$$E_{\text{aug}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n; \mathbf{w}) - y_n)^2 + \frac{\lambda}{N} \sum_{\ell, i, j} (w_{ij}^{(\ell)})^2$$

$$\frac{\partial E_{\text{aug}}(\mathbf{w})}{\partial W^{(\ell)}} = \underbrace{\frac{\partial E_{\text{in}}(\mathbf{w})}{\partial W^{(\ell)}}}_{\substack{\uparrow \\ \text{backpropagation}}} + \frac{2\lambda}{N} W^{(\ell)}$$

## Aspectos prácticos en entrenamiento

Early stopping: agregar un set de validación para monitorear durante el entrenamiento el error fuera de muestra.

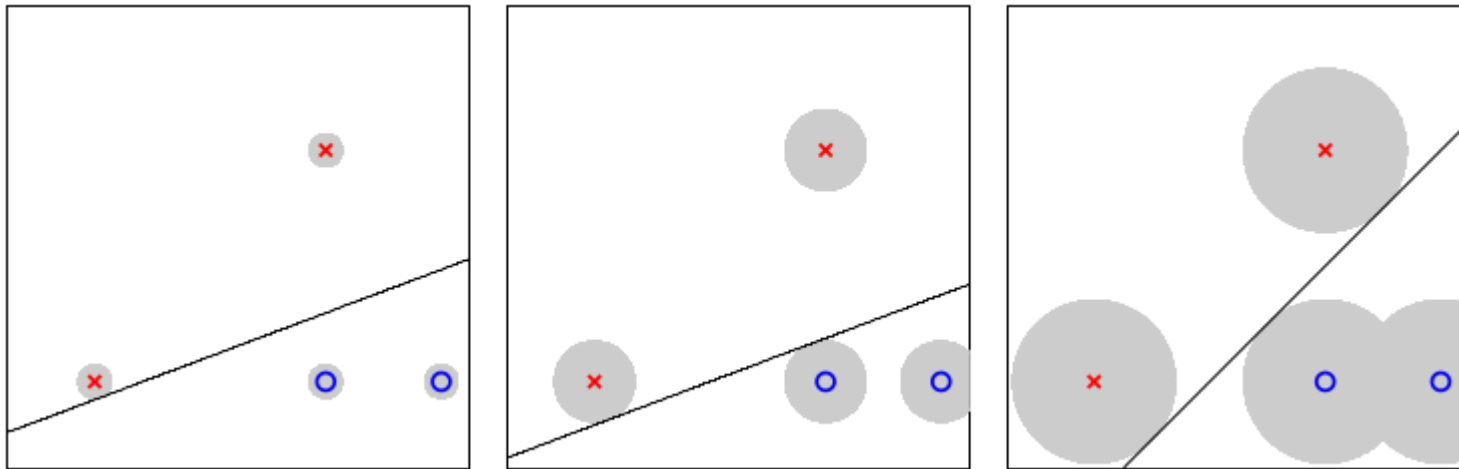


# Learning with kernels

## Separadores y datos ruidosos

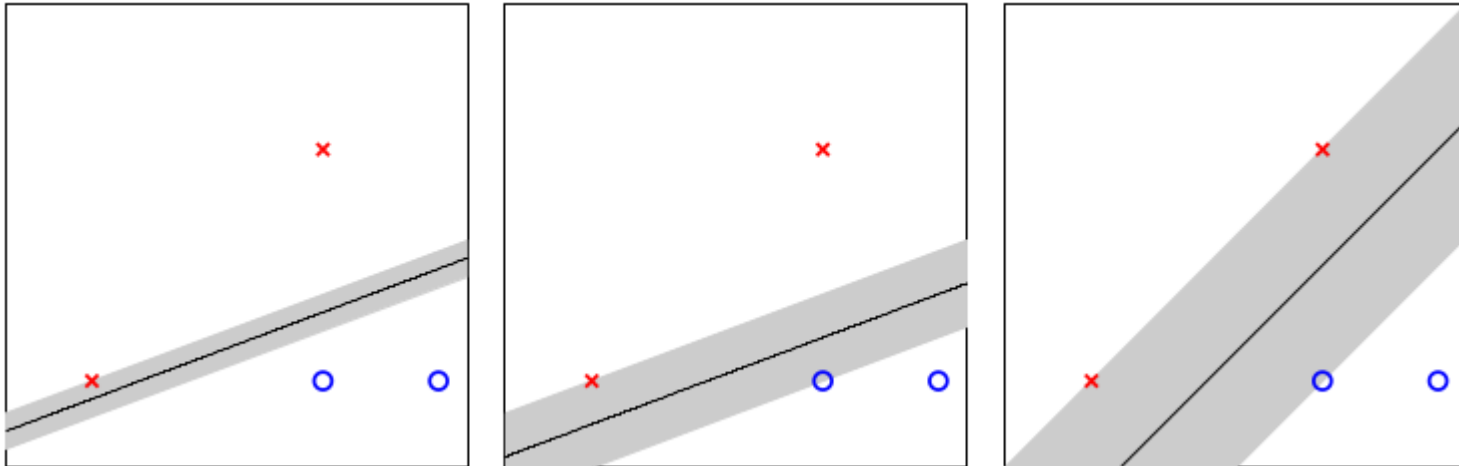
- Los datos pueden ser ruidosos (ruido de medición).
- Nuestros modelos debieran ser robustos a datos ruidosos.

Idea: La robustez al ruido tiene relación con considerar un margen de error para las mediciones.



## Separadores y datos ruidosos

Una idea análoga a márgenes para datos consiste en trabajar con **hiperplanos gruesos**, agregando un margen al separador.



# Separadores y datos ruidosos

Para trabajar con un hiperplano grueso, podemos usar el sesgo de una manera ingeniosa.

## Hiperplano estrecho

$$\mathbf{x} \in \{1\} \times \mathbb{R}^d; \mathbf{w} \in \mathbb{R}^{d+1}$$

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}; \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}.$$

$$\text{signal} = \mathbf{w}^T \mathbf{x}$$



El sesgo se codifica como una dimensión más

## Hiperplano grueso

$$\mathbf{x} \in \mathbb{R}^d; b \in \mathbb{R}, \mathbf{w} \in \mathbb{R}^d$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}; \quad \mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}.$$

bias  $b$

$$\text{signal} = \mathbf{w}^T \mathbf{x} + b$$



El sesgo aditivo interviene en el espacio de representación

# Separadores y datos ruidosos

## Hiperplano grueso

$$\mathbf{x} \in \mathbb{R}^d; b \in \mathbb{R}, \mathbf{w} \in \mathbb{R}^d$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}; \quad \mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}.$$

bias  $b$

$$\text{signal} = \mathbf{w}^T \mathbf{x} + b$$

↑  
sesgo

