

Teoría de generalización  $\rightarrow$  ¿Qué tan buena es la red al predecir?

¿Cuál es la influencia de la cantidad de ejemplos de entrenamientos?

$\bar{w}$  : conjunto de parámetros de la red

$\tilde{f}(\vec{x})$  : función blanco

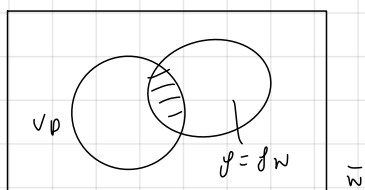
$\vec{x}$  : vectores de entrada

$(\vec{x}^u, y^u = \tilde{f}(\vec{x}^u)) \quad u=1, \dots, P$  pares de aprenden  $\rightarrow$  conjunto de entrenamientos

$\vec{x}^u$  son extraídos de  $P(\vec{x})$  (independientes)

distribución de proba de los parámetros

$f_{\bar{w}}(\vec{x})$  función implementada por  $\bar{w}$



$$V_0 = \int d\bar{w} P(\bar{w})$$

ninguna  
total  
en el  
espacio  
de aprendizaje

$$\Theta_f(\bar{w}) = \begin{cases} 1 & \text{si } f(\vec{x}) = f_w(\vec{x}) \\ 0 & \text{c.c.} \end{cases}$$

$$V(f) = \int d\bar{w} P(\bar{w}) \Theta_f(\bar{w})$$

$$\pm(f_w, \vec{x}^u) = \begin{cases} 1 & \text{si } f_w(\vec{x}) y^u \equiv \tilde{f}(\vec{x}^u) \\ 0 & \text{c.c.} \end{cases} \quad V_P = \int d\bar{w} P(\bar{w}) \left[ \prod_{u=1}^P I(f_{\bar{w}}, \vec{x}^u) \right]$$

$$V(P, f) = \int d\bar{w} P(\bar{w}) \Theta_f(\bar{w}) \left[ \prod_{u=1}^P \pm(f_{\bar{w}}, \vec{x}^u) \right]$$

$$V_P(f) = \left[ \int d\bar{w} P(\bar{w}) \Theta_f(\bar{w}) \right] \left[ \prod_{u=1}^P I(f, \vec{x}^u) \right]$$

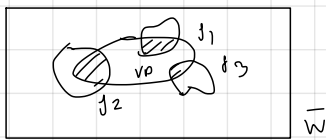
$\langle \dots \rangle$  sobre  $P(\vec{x})$

$$\langle V_P(f) \rangle = \left[ \int d\bar{w} P(\bar{w}) \Theta_f(\bar{w}) \right] \langle \prod_{u=1}^P I(f, \vec{x}^u) \rangle$$

$$= \left[ \int d\bar{w} P(\bar{w}) \Theta_f(\bar{w}) \right] \underbrace{\langle I(f, \vec{x}) \rangle^P}_{g(f)}$$

$$I(f, \vec{x}) = \begin{cases} 1 & \text{si } f(\vec{x}) = \tilde{f}(\vec{x}) \\ 0 & \text{c.c.} \end{cases}$$

$$\langle V(P, f) \rangle = V_0(f) g(f)^P$$



$$P_p(y) \propto v_p(y) g(y)^p$$

$$P_p(g) = \sum_j P_p(y) \delta(g - g(y)) \propto \sum_j v_p(y) g(y)^p \delta(g - g(y))$$

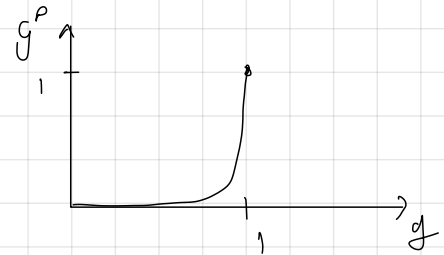
$$\propto g^p \underbrace{\sum_j v_p(y) \delta(g - g(y))}_{g_0(g)}$$

$$P_p(g) = \frac{g^p P_0(g)}{\int_0^1 g'^p P_0(g') dg'}$$

$$\langle g \rangle = \frac{\int_0^1 g P_p(g) dg}{\int_0^1 g^p P_0(g) dg}$$

$$= \frac{\int_0^1 g^{p+1} P_0(g) dg}{\int_0^1 g^p P_0(g) dg}$$

$$\approx \frac{\int_0^1 g^{p+1} dg \cancel{P_0(1)}}{\int_0^1 g^p dg \cancel{P_0(1)}} = \frac{1}{\frac{p+2}{p+1}} = \frac{1}{p+1}$$



$$= \frac{p+1}{p+2} = \frac{1+p}{1+2p} = \left(1 + \frac{1}{p}\right) \left(1 - \frac{2}{p}\right)$$

$$\approx 1 - \frac{1}{p}$$

es decir de  
como se tiene  
que incrementar  
el n° de ejemplos  
para mejorar la  
generalización

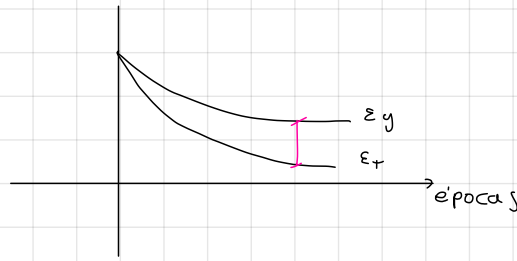
Interesa cual es la probabilidad del peor caso  
COTA DEL PEOR CASO  
Por simplicidad tomemos funciones booleanas

$\tilde{f}(\vec{x})$ : función booleana

$$(\vec{x}, \vec{y}^m \equiv \tilde{f}(\vec{x}^m) \quad m=1, \dots, p$$

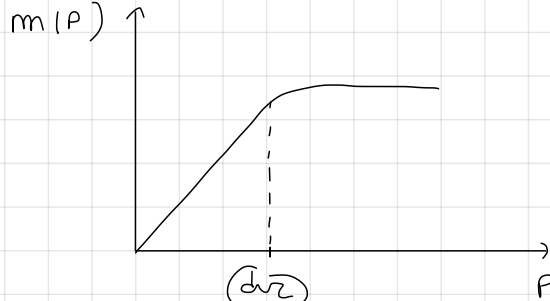
$$g(y) = \langle I(f, \vec{x}) \rangle$$

$g_P(f) = \left\langle \prod_{i=1}^P I(f, \tilde{x}^i) \right\rangle \rightarrow$  performance sobre el conjunto de entrenamientos



típicamente  
 $g_P(f) > g(f)$

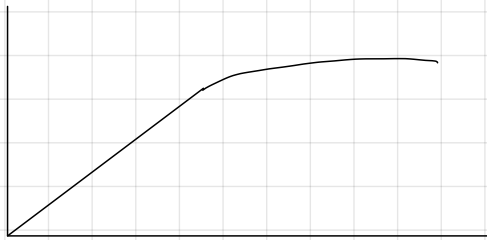
$$\text{PROB} \left( \max_f |g_P(f) - g(f)| \right) \leq 4 m(2, 0) e^{-\epsilon^2 P / 8} \rightarrow 0 \text{ as } P \rightarrow \infty \forall \epsilon$$



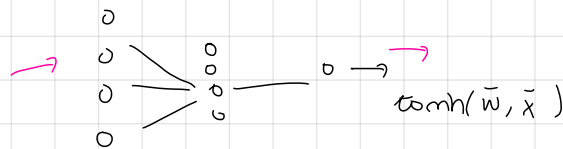
$$m(P) \leq P^{d_{vc}} + 1$$

$d_{vc}$   
 resumen de la complejidad de la red  
 potencia mínima de P para que exista

$\rightarrow$  dimensión Vapnik - Chernovskiz



perceptrón 2 capas  
 $2N = d_{vc}$   
 $d_{vc} = N m \lg(m)$

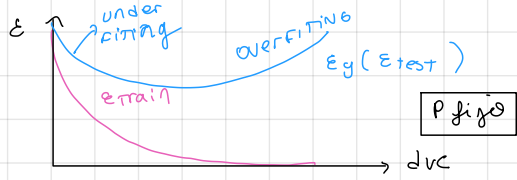


escritor contu  
 ver películas  
 cocinar  
 pedir verduras  
 elegir películas

$d_{vc}$  debe depender del tamaño de los pesos

1 capa continua  $d_{vc} \leq \min \{ D^2 / w^2, 2N + 1 \}$

- in psicología
- reunión matemática



Combiarmos la complejidad de la red

$\epsilon$  para tener una idea de poder decir  
crece

