

¿Porqué es difícil entrenar redes profundas?

Ahora tenemos una idea bastante ecabada de como funciona una red neuronal de un número arbitrario de capas con neuronas lineales o no lineales.

La esencia del método, como vimos, incluye

- disponer de un conjunto de datos previamente etiquetados con la respuesta correcta. Este conjunto de datos nos permitirá **entrenar** nuestra red.
- reservar al menos una parte de este conjunto de datos etiquetados que usamos para evaluar como resuelve la red los casos que nunca he visto. Solo si hace bien este caso, diremos que la red ha **aprendido** y es capaz de **generalizar**.
- Debemos definir una función error $E(\{\bar{w}\})$ que depende de todos los espalamientos sinápticos de nuestra red. Por ahora asumimos un error cuadrático.

- Preparamos la red eligiendo aleatoriamente todos los pesos sinápticos.
- Ahora aplicamos el descenso por el gradiente determinista. Le presentamos secuencialmente los inputs del conjunto de entrenamiento y obtenemos los resultados W_i^H ($i=1,2,\dots,M$). Con estos calculamos E , con E calculamos los incrementos de cada peso ΔW_{pq}

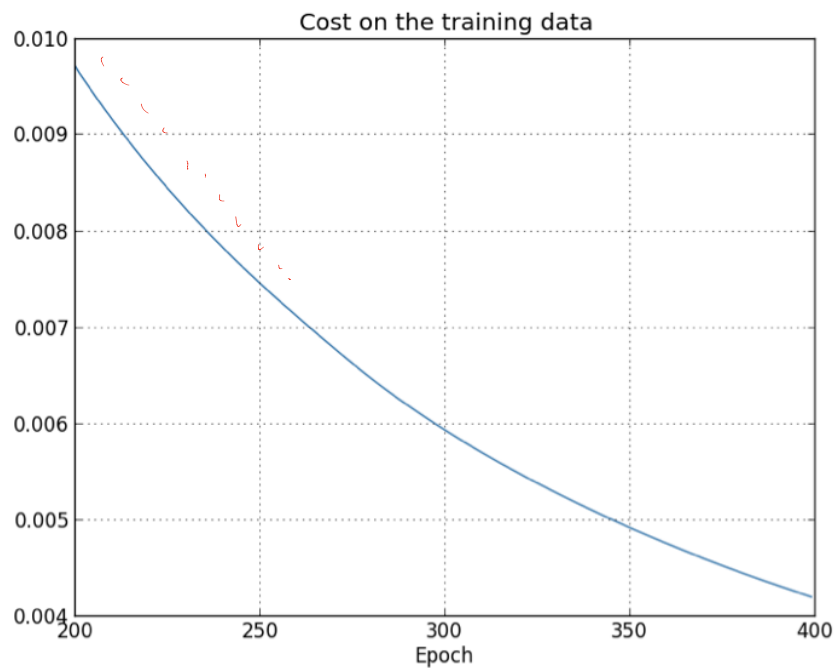
$$\Delta W_{pq} = -\eta \frac{\partial E}{\partial W_{pq}}$$

Cuando le hemos presentado los p patrones del conjunto de entrenamiento, hemos pasado **1 ÉPOCA**. Calculamos el valor E en la época 1 .

- Repetimos el último paso incrementando la época en 1

$$\text{época} =: \text{época} + 1$$

y monitoreamos el valor de E en función de **época**. Sabemos que este valor debe decaer o mantenerse constante.



Ejemplo: E debe decaer o mantenerse constante

No detenemos cuando el error alcance cierta tolerancia preestablecida o cuando alcancemos cierto número de épocas.

Ahora hacemos una evaluación de la capacidad de aprendizaje, época por época o al final.

Se puede mostrar que con una capa oculta con un número arbitrario de neuronas es posible aproximar tan bien como queramos una función dada, bajo condiciones específicas. Pero esto requiere un número exponencialmente grande de neuronas cuando aumenta la complejidad del problema.

En una red de una capa oculta vemos que la salida O_i depende de la entrada a través de las neuronas de la capa intermedia V

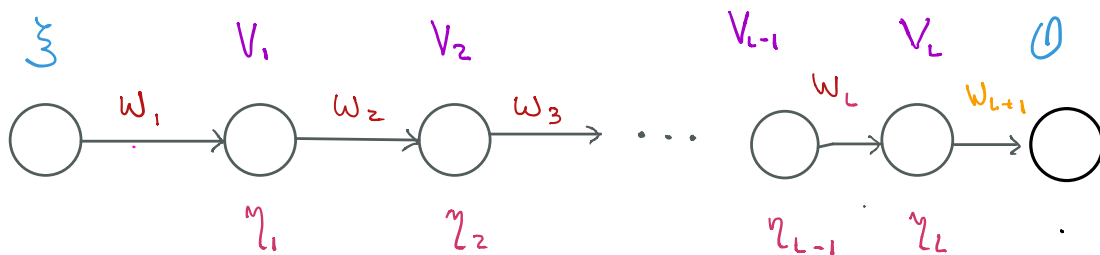
$$O_i = g \left(\sum_j W_{ij} \cdot g \left(\sum_{k=1} W_{jk} \xi_k \right) \right)$$

Incluso podemos tener funciones diferentes en la capa intermedia y en la de salida:

$$O_i^h = g_1 \left(\sum_j W_{ij} \cdot \overbrace{g_2 \left(\sum_{k=1} W_{jk} \xi_k^h \right)}^{V_j(h_j)} \right)$$

Con estos O_i^k calculamos los δ^k y con ellos "volvemos" de la salida a la entrada cambiando los complementos.

La red más simple, y por cierto inútil, es una red de muchos capas con una única neurona por capa



$$\begin{aligned}
 O &= g_{L+1}(w_{L+1}V_L - \eta_L) = g_{L+1}(w_{L+1}g_L(w_LV_{L-1} - \eta_{L-1}) - \eta_L) \\
 &= g_{L+1}(w_{L+1}g_L(w_Lg_{L-1}(w_{L-1}V_{L-2} - \eta_{L-2}) - \eta_{L-1}) - \eta_L) \\
 &= g_{L+1}(w_{L+1}g_L(w_Lg_{L-1}(w_{L-1}g_{L-2}(w_{L-2}g_{L-3}(w_{L-3}g_{L-4} \dots) - \eta_L) - \eta_{L-1}) - \eta_L)
 \end{aligned}$$

Hacia adente componemos funciones

$$\frac{\partial E}{\partial w_1} = g'_1(h_1) \times w_2 \times g'_2(h_2) \times \dots \times w_L g'_L(h_L) \times \underbrace{[\delta^k - O^k]}_{\text{diferencia entre deseado y real}}$$

Hacia atrás multiplicamos derivadas de funciones

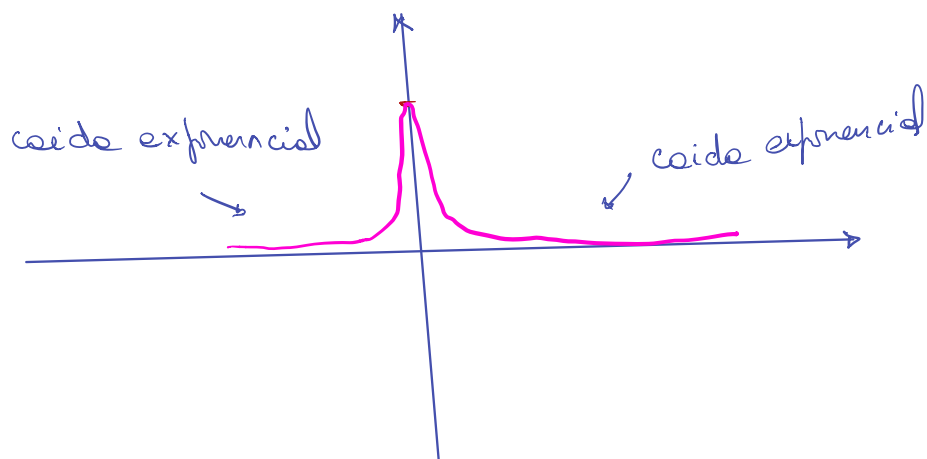
Miremos un segundo las dos funciones que se usaban originalmente:

$$g(h) = \tanh(h) = \frac{e^h - e^{-h}}{e^h + e^{-h}} = \frac{\operatorname{senh}(h)}{\operatorname{cosh}(h)}$$

$$g'(h) = \frac{(e^h + e^{-h})(e^h - e^{-h}) - (e^h - e^{-h})(e^h + e^{-h})}{(e^h + e^{-h})^2}$$

$$= \frac{(e^h - e^{-h})^2 - (e^h + e^{-h})^2}{(e^h + e^{-h})^2}$$

$$= 1 - \tanh^2(h) = \underline{1 - g^2(h)}$$



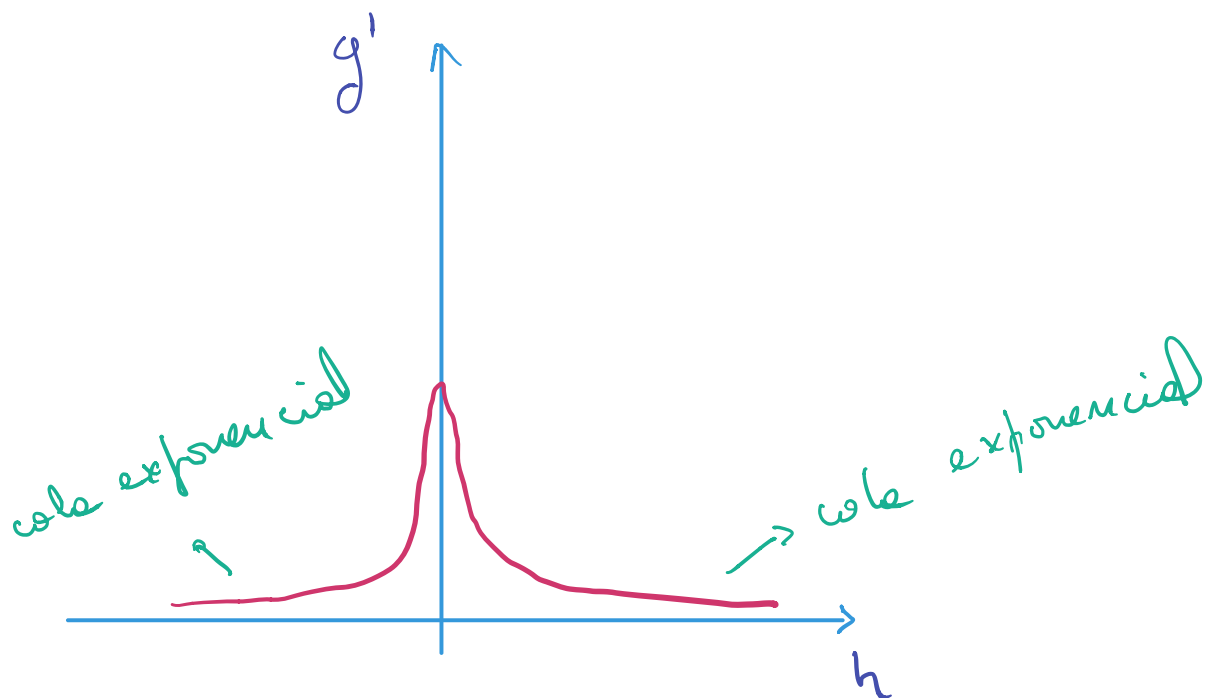
Para neuronas sigmoideas:

$$g(h) = \frac{1}{1+e^{-h}}$$

$$g'(h) = \frac{-(-e^{-h})}{(1+e^{-h})^2} = \frac{e^{-h}}{(1+e^{-h})^2}$$

$$= \frac{1}{(1+e^{-h})} \frac{e^{-h}}{(1+e^{-h})} = g(h) \cdot \left(1 - \frac{1}{1+e^{-h}}\right)$$

$$= g(h) (1 - g(h))$$



Noten que $g'(u)$ tiene una imagen exponencialmente pequeña salvo en un entorno de cero con un ancho típico de uno. Entonces cuando obtenemos muchas neuronas en muchas capas, a medida que avanza el back-propagation hacia atrás, los incrementos se hacen exponencialmente pequeños y la red no puede aprender. Esto se llama

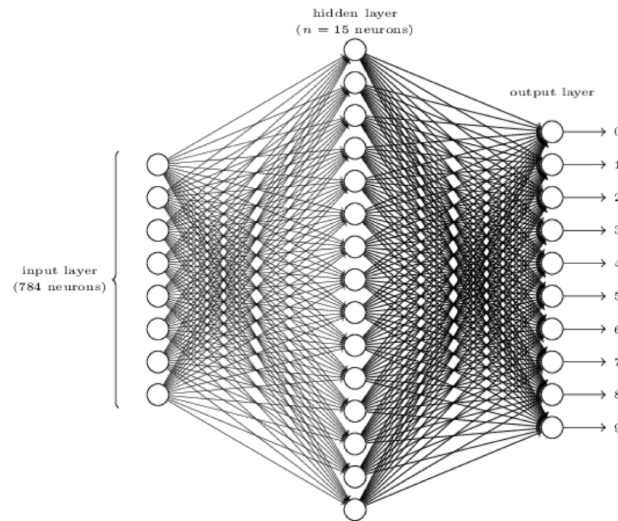
SUPRESIÓN DEL GRADIENTE

Miemos un ejemplo.

Se hizo una red multicapa para aprender a distinguir dígitos escritos a mano sacados de la base de datos MNIST.

La entrada son 724 neuronas correspondientes a los 28×28 píxeles de las fotos de MNIST por medio de matriz bidimensional a vector (lineal).

Hacemos redes feed-forward cada vez más profundas para aprender a reconocer los 10 dígitos (MNIST).

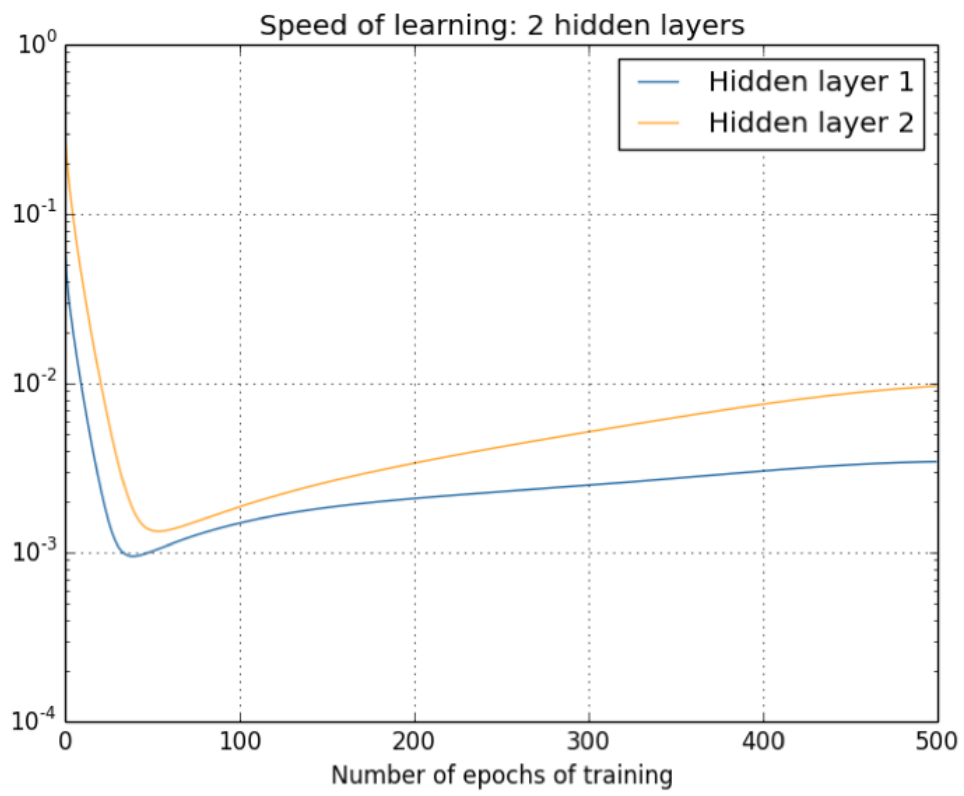


Miremos qué sucede con un acoplamiento en particular a medida que hacemos más profunda la red, comenzando con 2 capas ocultas.

Miremos como cambia la región de cambio del umbral de la primera neurona de la primera capa oculta ($\gamma = w_0^{(1)}$)

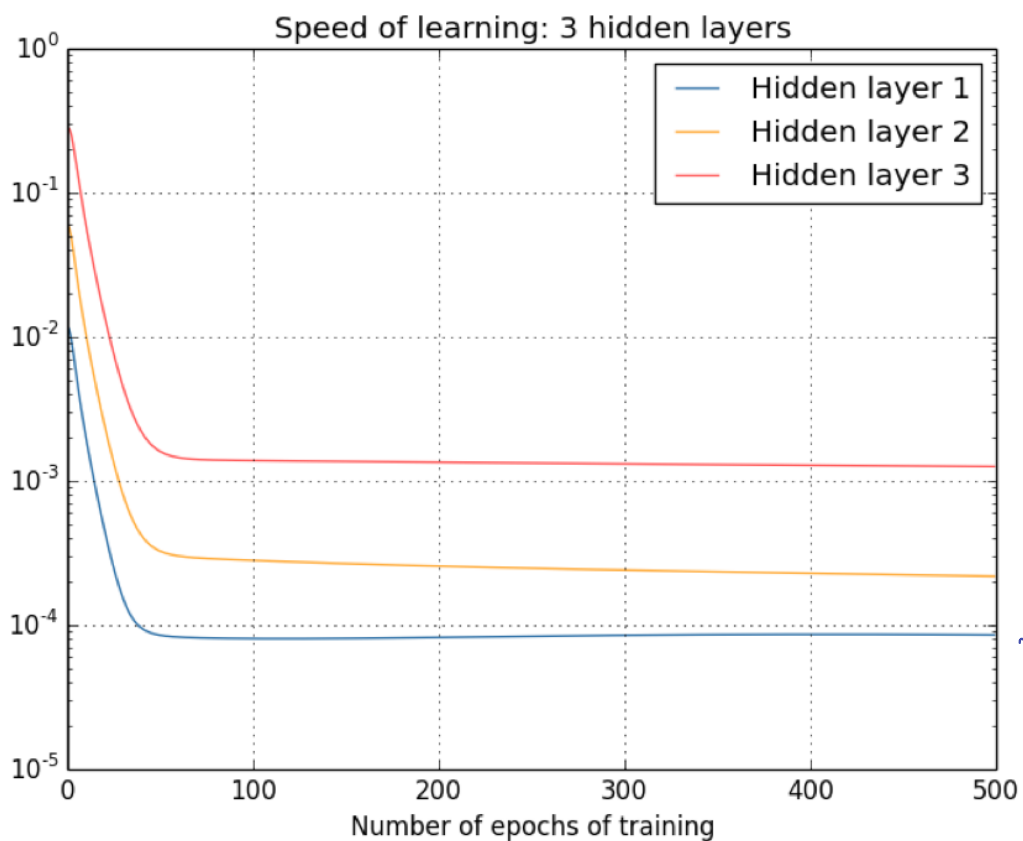
Esta velocidad se representa con los curvas azules
Cada capa oculta tiene 100 neuronas.

2 CAPAS OCULTAS

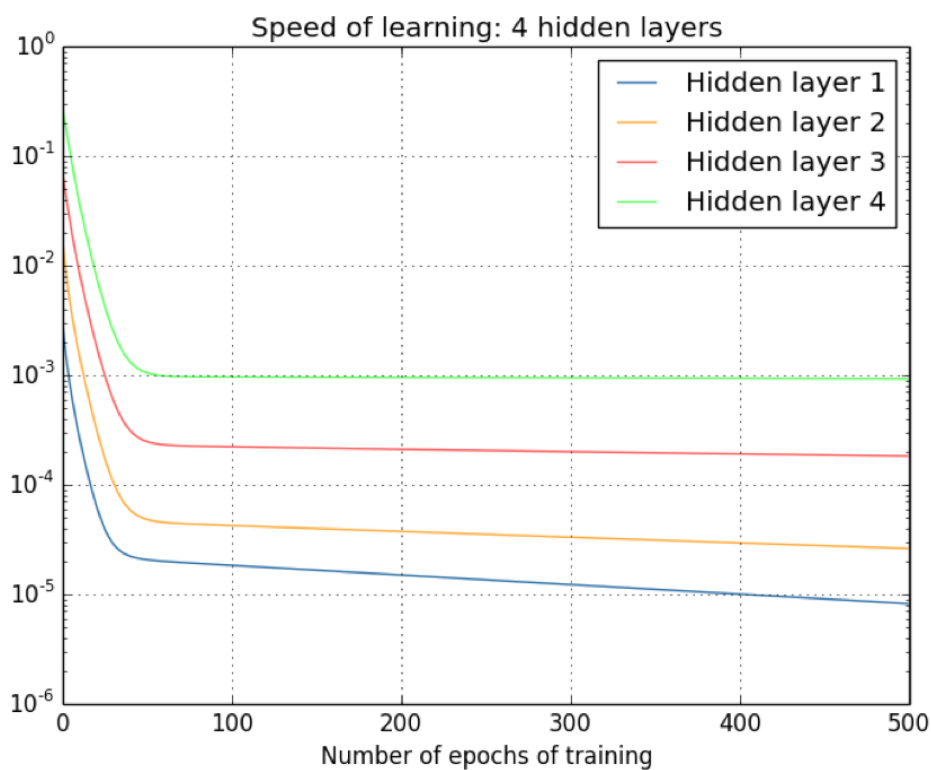


$\sim 2,6 \times 10^{-2}$

3 CAPAS OCULTAS



4 CAPAS OCULTAS



Vemos que cuanto más capas ponemos, la última capa de sinapsis aprende siempre a la misma velocidad (aproximadamente) y la primera cada vez más lentamente.

SUPRESIÓN DEL GRADIENTE

Esto es culpa de varias cosas:

- * las funciones sigmoidea y tangente tienen derivadas pequeñas, y cada vez que vamos hacia atrás, la razón de cambio se hace multiplicando más derivadas.
- * el método de gradiente queda preso en mínimos locales
- * la razón de aprendizaje es fija.
- * los W iniciales son distribuidos en un hiperespacio de muchas dimensiones

En la próxima clase veremos como solucionar estos problemas.