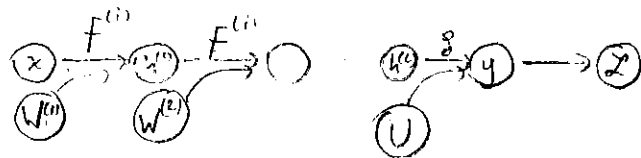
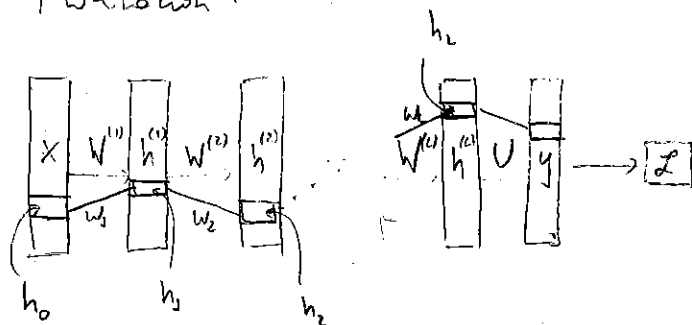


Optimización (Iniciolización)

- Idea: mejorar los pesos de entre 10 minutos, volver + Iniciolización, normalización
- + Algoritmo para descenso de gradiente no efectivo

Motivación:



sea h_i una unidad cualquiera de la capa encubierta de $h^{(i)}$

$$h_i = f_i(h_{i-1}w_i + b_i) \quad \text{para } i \in \{1, \dots, L\}$$

algo que no depende de w_i ni de h_{i-1}

si hago Back Propagation desde z hasta w_1 por la rama identificada

$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial h_L} \cdot \frac{\partial h_L}{\partial h_{L-1}} \cdot \frac{\partial h_{L-1}}{\partial h_{L-2}} \cdot \frac{\partial h_{L-2}}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial h_L} \left(\prod_{j=2}^L \frac{\partial h_j}{\partial h_{j-1}} \right) \cdot \frac{\partial h_1}{\partial w_1}$$

$$\frac{\partial h_j}{\partial h_{j-1}} \Rightarrow \frac{\partial f_j(h_{j-1}w_j + b_j)}{\partial h_{j-1}} = f'_j(h_{j-1}w_j + b_j) \cdot w_j$$

llamemos a esto u_j

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial h_L} \left(\prod_{j=2}^L f'_j(u_j) \right) \left(\prod_{j=2}^L w_j \right) \frac{\partial h_1}{\partial w_1}$$

$$h_1 = f_1(xw_1 + b_1)$$

$$\frac{\partial h_1}{\partial w_1} = f'_1(u_1) x$$

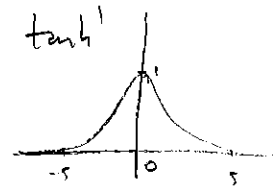
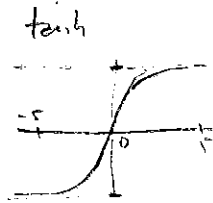
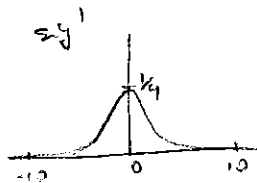
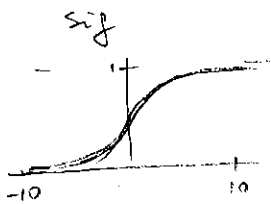
$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial h_L} \left(\prod_{j=1}^L f'_j(u_j) \right) \left(\prod_{j=2}^L w_j \right) x$$

- basta que uno de estos valores sea 0 para que el gradiente

sea 0 $\Rightarrow \frac{\partial \mathcal{L}}{\partial w_1} = 0 \Rightarrow$ no hay aprendizaje para w_1

- si varios de ellos son valores pequeños $\Rightarrow \frac{\partial \mathcal{L}}{\partial w_1}$ será muy (muy) pequeño si L es muy grande

coser



$$\frac{\partial \mathcal{L}}{\partial w_j} = \frac{\partial \mathcal{L}}{\partial h_L} \cdot \prod_{j=1}^L f'(u_j) \cdot \prod_{j=2}^L w_j \cdot x \quad \text{sig}'(u_j) \leq \frac{1}{4}$$

② sig. que nunca, sigmoid

\Rightarrow si los u_j 's son muy grandes

el aprendizaje ni siquiera comenzará
 \therefore mucho cuidado al inicializar!!

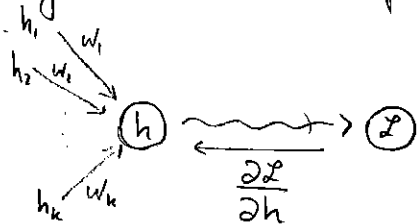
① \Rightarrow si los w_j son todos muy pequeños
 $\Rightarrow \prod_{j=2}^L w_j$ será muy (muy) pequeño

\Rightarrow hay que tratar de mantenerse en el área activa del sigmoid (cerca del 0)

③ sig. sigmoid todavía sin importar el valor de $u_j \Rightarrow \text{sig}'(u_j) \leq \frac{1}{4}$

$\Rightarrow \prod_{j=1}^L \text{sig}'(u_j) \leq \left(\frac{1}{4}\right)^L \Rightarrow$ si L es muy grande el gradiente se irá a 0 (vanishing gradient problem). Por lo mismo para todos los caminos que lleven a w_1

④ sigmoid tiene otro problema más:



$$h = \text{sig}(h_1 w_1 + h_2 w_2 + \dots + h_k w_k + b)$$

$$\frac{\partial \mathcal{L}}{\partial w_i} = h_i \cdot \text{sig}'(u) \cdot \frac{\partial \mathcal{L}}{\partial h}$$

si todos los h_i 's son positivos \Rightarrow

el signo de $\frac{\partial \mathcal{L}}{\partial w_i}$ dependerá solo de $\text{sig}'(u) \cdot \frac{\partial \mathcal{L}}{\partial h}$

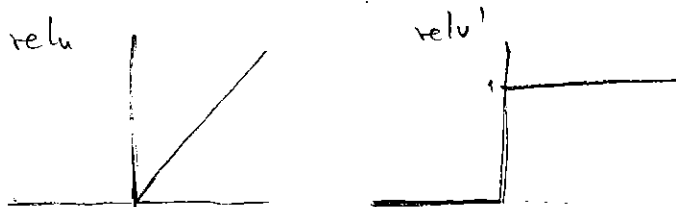
\Rightarrow todos los $\frac{\partial \mathcal{L}}{\partial w_i}$ tendrán siempre el mismo signo.

Note que $\text{sig}(u)$ es siempre positivo \Rightarrow en cada iteración de una red con activación sig todos los pesos asociados a una unidad tendrán gradiente con el mismo signo. (En general no gustaría que para cada unidad, los inputs estén alrededor en el 0!!!)

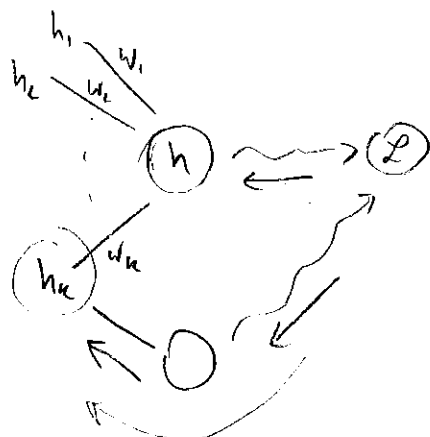
★ Morelijo 1: NUNCA USAR sigmoid EN DEEP LEARNING

Para tanh \Rightarrow tenemos los problemas ② y ③ pero no tenemos el ④

★ Morelijo 2: entre tanh y sigmoid, siempre preferir tanh



con relu no hay derivados en 0
de gradiente $\frac{\partial \mathcal{L}}{\partial w_i}$ no tiende a 0
en un lado mas "directo"



$$\frac{\partial \mathcal{L}}{\partial w_i} = \frac{\partial \mathcal{L}}{\partial h} \cdot \underbrace{\text{relu}'(u)}_{0 \text{ o } 1} \cdot h_i$$

- si es 0 $\Rightarrow w_i$ no se actualizara pero no significa que nunca mas se actualizara porque las h_i 's pueden seguir recibiendo gradiente por otros caminos

gradiente que puede cambiar el valor de activación de la relu

Problema: relu tiene el problema (4) de sigmoid (los inputs de cada unidad no están centrados en 0).

En la práctica relu funciona mejor que tanh !!! (de hecho es uno de las razones del desempeño de redes neurales profundas)

★ Consejo 3: SIEMPRE USAR RELU EN DEEP LEARNING

hay otra familia a relu que centran su output en 0:

- Leaky RELU
- ELU, CELU etc.

Observación: los pesos iniciales pueden impactar sustancialmente en el entrenamiento

recuerda: la media y la varianza de un conjunto

+ Inicialización:

los cálculos siempre se vera como $f(h_1 w_1 + h_2 w_2 + \dots + h_k w_k)$
queremos que $z = h_1 w_1 + h_2 w_2 + \dots + h_k w_k$ no sea demasiado variable a pesar de que los elegimos al azar.

En notación matricial más

$$W^{(i)} = \text{randn}(d_{i-1}, d_i) \cdot \frac{1}{\sqrt{d_{i-1}}}$$

intento: $w_i = \underbrace{\text{randn}}_{\text{distribución normal}} \cdot \sqrt{\frac{1}{k}}$

Inicialización de Xavier

asegure que la media es 0 y la varianza es $\frac{1}{k}$
 \Rightarrow varianza de la suma es 1
 $k =$ dimensión de entrada a la capa

- Xavier funciona bien en general, pero empíricamente se ha visto que para relu una mejor inicialización es

$$W_i = \text{randnormal} \cdot \sqrt{\frac{2}{k}} \quad \left[\text{En notación } W^{(i)} = \text{randn}(d_{i-1}, d_i) * \sqrt{\frac{2}{d_{i-1}}} \right]$$

- Otra propuesta incluye dividir por la suma de las dimensiones

$$W^{(i)} = \text{randn}(d_{i-1}, d_i) * \sqrt{\frac{1}{d_{i-1} + d_i}}$$

- Todo esto es empírico, también se podría "buscar" la varianza correcta como un hiperparámetro

En general queremos que los outputs de nuestras capas sean 0-centered (también nos gustaría que tuvieran la misma varianza en todas las direcciones). La primera capa es el input y ahí podemos argumentar !!!

$$D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\} \text{ y nos } X \text{ tenemos con los } x^{(i)} \text{ } X(N, C)$$

$$\mu := \frac{1}{N} \sum_{i=1}^N x^{(i)} \quad \left(= \frac{1}{N} \text{sum}_1(X) \right) \quad \left(\downarrow \begin{matrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(N)} \end{matrix} \right) \leftarrow \text{media por cada feature}$$

$$X := X - \mu \quad \leftarrow \text{resta la media componente a componente a todos los ejemplares}$$

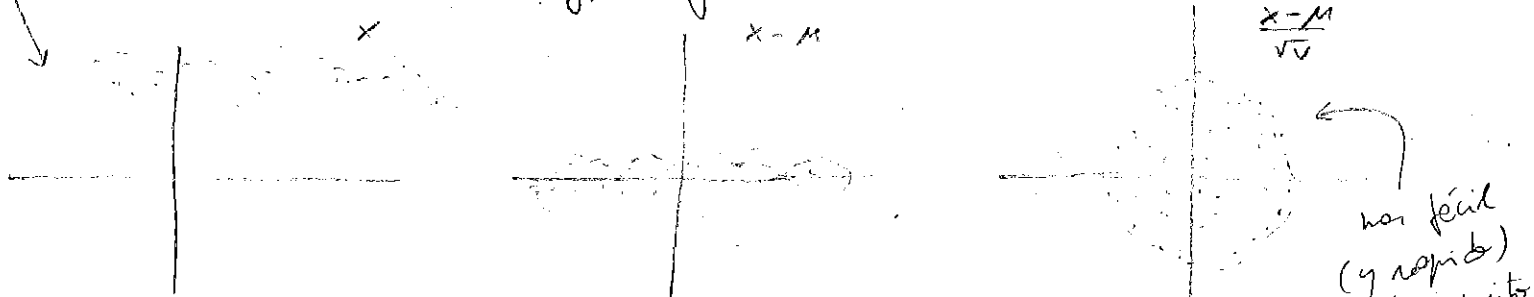
$$V := \frac{1}{N} \sum_{i=1}^N (x^{(i)} * x^{(i)}) \quad \left(= \frac{1}{N} \text{sum}_1(X * X) \right)$$

$$X := \frac{X}{\sqrt{V}} \quad \leftarrow \text{broadcast} \quad \left(= X * 1 / \text{sqrt}(V) \right)$$

$$\begin{aligned} & q_1, q_2, \dots, q_N \\ & \mu = \frac{1}{N} \sum q_i \\ & v = \frac{1}{N} \sum (q_i - \mu)^2 \\ & q_i' = \frac{q_i - \mu}{\sqrt{v}} \quad \left(\begin{array}{l} \text{normalizado} \\ \text{media } 0 \\ \text{varianza } 1 \end{array} \right) \end{aligned}$$

★ Consejo: siempre normalizar el input (no de no)

- en el caso de imágenes generalmente sólo se centra (valores entre 0-255)

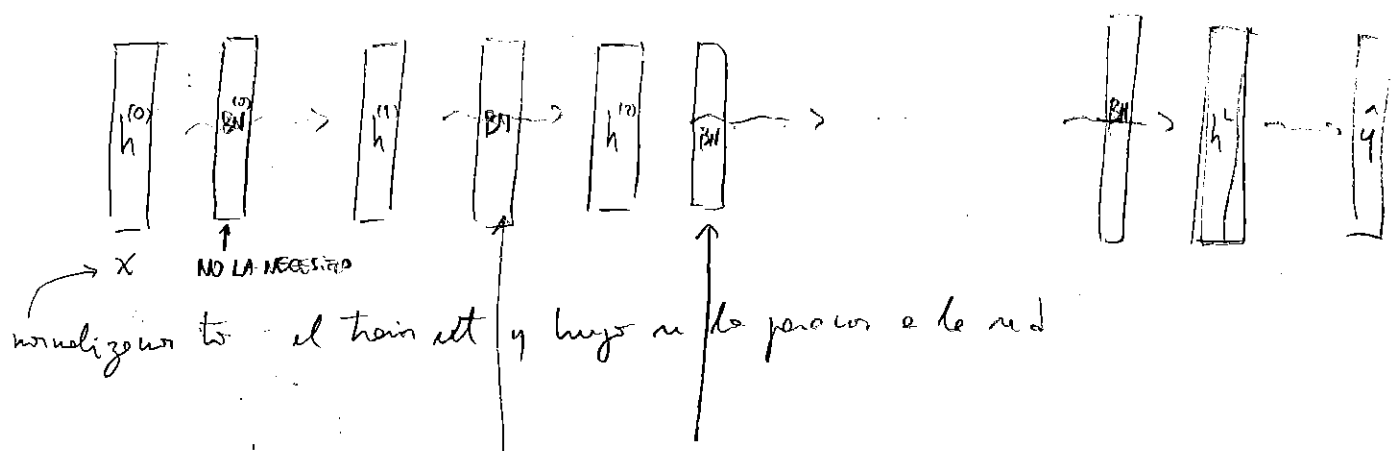


Nota: no normalizem todos los datos del training, el μ y v computados ahí se usan para normalizar en test time

no fácil
(y rápido)
el entrenamiento

Normalización por Paquetes (Batch Normalization)

Idea: normalizar el input tiene buenas propiedades para el entrenamiento. Podemos hacer algo similar dentro de la red?



Traer de normalizar entre media de los capas!!! ¿Cómo?

No olvidarnos que los ejemplos vienen en paquetes

$$h^{(i+1)} = f^{(i+1)} \left(\underbrace{h^{(i)}}_{\substack{\uparrow \\ \text{input de la capa } i+1}} W^{(i+1)} + b^{(i+1)} \right)$$

input de la capa $i+1 \Rightarrow$ normalizar $h^{(i)}$ a lo largo del batch

$$h^{(i)}(B, d_i) \leftarrow \text{dimensiones}$$

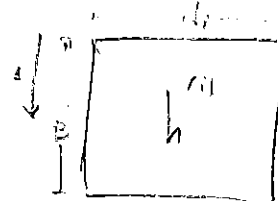
$$\mu^{(i)} = \frac{1}{B} \sum_{j=1}^B h_{j,:}^{(i)} = \frac{1}{B} * \text{sum}_{\text{axis } 0} (h^{(i)})$$

$$\mu^{(i)}(d_i) \leftarrow \text{dimensiones}$$

$$v^{(i)}(d_i) \leftarrow \text{dimensiones}$$

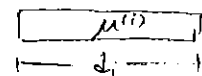
$$v^{(i)} = \frac{1}{B} \sum_{j=1}^B \left(\left(h_{j,:}^{(i)} - \mu^{(i)} \right)^2 \right)$$

↑
suma de vectores
broadcast!!!



$$\tilde{h}^{(i)} = \frac{h^{(i)} - \mu^{(i)}}{\sqrt{v^{(i)} + \epsilon}}$$

← normalice los elementos del batch
broadcast!! (ej: normalización por cada feature)



Quiero 3) $h^{(i+1)} = f^{(i+1)} \left(\tilde{h}^{(i)} W^{(i+1)} + b^{(i+1)} \right)$

Es derivable! puede seguir usando Back Prop

$$\frac{\partial h^{(i+1)}}{\partial h^{(i)}} = \frac{\partial h^{(i+1)}}{\partial \tilde{h}^{(i)}} \cdot \frac{\partial \tilde{h}^{(i)}}{\partial h^{(i)}}$$

Posible problema: entrenar logrando la capacidad de la red!!

$\tilde{h}^{(i)}$ como transformación es lineal
 \Rightarrow mi red lo podría haber aprendido

Solución: dándole la capacidad a la red \Rightarrow permitirle que
 desbaje la normalización si así lo quiere

$$\hat{h}^{(i)} := \gamma^{(i)} * \tilde{h}^{(i)} + \beta^{(i)}$$

↓
 nuevo parámetro de la red!!!

$$\begin{matrix} \tilde{h}^{(i)}(B, d_i) \\ \gamma^{(i)}(d_i) \\ \beta^{(i)}(d_i) \end{matrix} \leftarrow \begin{matrix} \text{dimension} \end{matrix}$$

Se actualizan igual que todos los otros calculando primero el gradiente

Nueva fórmula

$$h^{(i)} = f^{(i)}(\hat{h}^{(i-1)} W^{(i)} + b^{(i)})$$

$$\hat{h}^{(i)} = \gamma^{(i)} \text{norm}(h^{(i)}) + \beta^{(i)}$$

$$\text{norm}(h^{(i)}) := \frac{h^{(i)} - \mu^{(i)}}{\sqrt{v^{(i)} + \epsilon}}$$

↑ estabilidad numérica

Observaciones:

- Paper original (2015) propone normalización antes de la no linealidad

$$u^{(i)} = h^{(i-1)} W^{(i)} + b^{(i)}$$

$$\hat{u}^{(i)} = \gamma^{(i)} \text{norm}(u^{(i)}) + \beta^{(i)}$$

$$h^{(i)} = f^{(i)}(\hat{u}^{(i)})$$

← en este caso $b^{(i)}$ y $\beta^{(i)}$ funcionan como términos de bias redundantes pueden obviarse de los $b^{(i)}$'s (solo usar $\beta^{(i)}$'s)

- En test time puede no tener sentido usar $\mu^{(i)}$ y $v^{(i)}$ calculados desde un batch \Rightarrow debo fijar esos valores al finalizar el entrenamiento. Una forma es ir calculando un promedio de $\mu^{(i)}$, $v^{(i)}$ durante el entrenamiento, digamos $\bar{\mu}^{(i)}$, $\bar{v}^{(i)}$, y usar ese valor en test time.

- La opción más usada: "promedio móvil exponencial"

$$\begin{aligned} \bar{\mu}^{(i)t} &= \delta \bar{\mu}^{(i)t-1} + (1-\delta) \mu^{(i)t} \\ \bar{v}^{(i)t} &= \delta \bar{v}^{(i)t-1} + (1-\delta) v^{(i)t} \end{aligned}$$

$0 \leq \delta \leq 1$ $\delta \rightarrow 1 \Rightarrow$ recuerda desde el principio
 $0 \leq \delta \leq \frac{1}{2} \Rightarrow$ me quedo con el más reciente

valor típico: $\delta = 0.9 \approx$ promedio ponderado de

($\delta = 0.5 \Rightarrow$ los 2, $\delta = 0.98 \Rightarrow$ últimos 50) opp. los últimos 10

- En cada iteración calculo $\mu^{(i)}$ actualizo $\bar{\mu}^{(i)}$ como
 $\bar{\mu}^{(i)} := \delta \bar{\mu}^{(i)} + (1-\delta) \mu^{(i)}$
 igual pero $\bar{v}^{(i)}$

Implementación de hipótesis