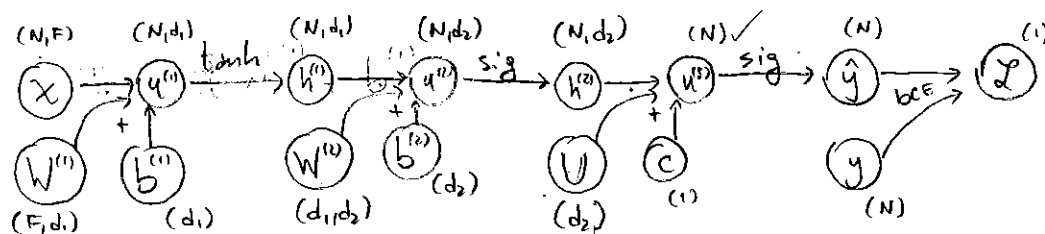


4) Ejemplo de Back Propagation completo a mano

①



$$u^{(1)} = x W^{(1)} + b^{(1)} \quad u^{(2)} = h^{(1)} W^{(2)} + b^{(2)} \quad u^{(3)} = h^{(2)} U + c$$

$$h^{(1)} = \tanh(u^{(1)}) \quad h^{(2)} = \text{sig}(u^{(2)}) \quad \hat{y} = \text{sig}(u^{(3)})$$

$$\text{error}(\hat{y}^{(i)}, y^{(i)}) = - (y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$$

$$L = \frac{1}{N} \sum_{i=1}^N \text{error}(\hat{y}^{(i)}, y^{(i)}) \quad H(P_{y^{(i)}}, P_{\hat{y}^{(i)}})$$

Interludio: Cross Entropy. Dadas dos distribuciones de probabilidad P, Q sobre un dominio finito D la entropía cruzada entre P y Q (en ese orden) está dada por

$$- \sum_{v \in D} P(v) \log(Q(v))$$

← idea: cuanto más "dependioso" si somos que más datos distribuyen como Q cuando efectivamente distribuyen como P

En teor. información la entropía de una (variable aleatoria con) distribución de probabilidad P se calcula como

$$\sum_{v \in D} P(v) \underbrace{I(v)}_{\text{información que entrega el observador } v} = \sum_{v \in D} P(v) \log\left(\frac{1}{P(v)}\right) = - \sum_{v \in D} P(v) \log(P(v))$$

Promedio de la información obtenida al observar un valor al azar de la distribución. $\Rightarrow - \sum_{v \in D} P(v) \log(Q(v))$ promedio de la información obtenida al observar un valor al azar según P si impongo (enormemente) que los valores distribuyen según Q

Cantidad promedio de bits por evento que se necesitan (como mínimo) para codificar una fuente de datos según P . $\Rightarrow - \sum_{v \in D} P(v) \log(Q(v))$ es una medida de la cantidad de bits que se desperdician si se impone (enormemente) que los datos distribuyen según Q

Volviendo a nuestro ejemplo: menor entropía cruzada con error (por cada ej.)

- imponemos que sabemos que el resultado de la red debe ser 1 ($y=1$) y nuestra red entrega el valor 0,89

\Rightarrow entropía cruzada entre el valor empírico y la predicción es

$$- 1 \cdot \log(0.89) = 0,117 \quad \leftarrow \text{error}$$

(y si nuestra red entrega el valor 0,32?)

$$- 1 \cdot \log(0.32) = 1,139 \quad \leftarrow \text{error}$$

- imponemos que sabemos que el resultado de la red debe ser 0 ($y=0$) y nuestra red entrega el valor 0,89

$$- 1 \cdot \log(1-0.89) = 2,207 \quad \leftarrow \text{error}$$

(y si entrega el valor 0,32?)

$$- 1 \cdot \log(1-0.32) = 0,386 \quad \leftarrow \text{error}$$

En general, si y es lo que esperamos ver ($y=0$ o 1) y la predicción de la red es \hat{y} (\hat{y} está entre 0,1)

$$\text{error}(\hat{y}, y) = \begin{cases} -\log(\hat{y}) & \text{si } y=1 \\ -\log(1-\hat{y}) & \text{si } y=0 \end{cases}$$

$$\text{error}(\hat{y}, y) = - (y \log(\hat{y}) + (1-y) \log(1-\hat{y}))$$

$$\begin{array}{ccccc} \uparrow & & \uparrow & & \uparrow \\ P_{\text{emp}}(y=1) & & P_{\text{emp}}(y=0) & & P_{\text{NN}}(y=0) \\ & \nearrow & & \nearrow & \\ & P_{\text{NN}}(y=1) & & & \end{array}$$

\Rightarrow nuestra función de error es la entropía cruzada entre la distribución de probabilidad "empírica" y la generada por la red (para cada ejemplo)

4) Ejemplo Back Prop. (continuación)

③

Neurite we $\frac{\partial \mathcal{L}}{\partial \sigma}$ para cada parámetro σ . Lo hacemos con programación dinámica sobre el grafo de computación de la red.

$$\frac{\partial \mathcal{L}}{\partial \hat{y}} \quad \mathcal{L} = -\frac{1}{N} \sum_{i=1}^N (y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)) \quad \begin{matrix} \hat{y}_i = \hat{y}^{(1)} \\ y_i = y^{(1)} \\ \uparrow \\ \text{vector} \end{matrix}$$

$$\frac{\partial \mathcal{L}}{\partial \hat{y}_i} = -\frac{1}{N} \left(\frac{y_i}{\hat{y}_i} + \frac{(1-y_i)}{(1-\hat{y}_i)} \cdot -1 \right) = \frac{1}{N} \left(\frac{(1-y_i)}{(1-\hat{y}_i)} - \frac{y_i}{\hat{y}_i} \right)$$

$$\frac{\partial \mathcal{L}}{\partial u^{(3)}_i} = \frac{\partial \mathcal{L}}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial u^{(3)}_i} \quad \hat{y}_i = \text{sig}(u^{(3)}_i) \Rightarrow \hat{y}_i = \text{sig}(u^{(3)}_i)$$

suma de Einstein !!!

$$= 0 \text{ si } j \neq i \Rightarrow \frac{\partial \mathcal{L}}{\partial u^{(3)}_i} = \frac{\partial \mathcal{L}}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial u^{(3)}_i}$$

$$\frac{\partial \hat{y}_i}{\partial u^{(3)}_i} = \frac{\partial \text{sig}(u^{(3)}_i)}{\partial u^{(3)}_i} = \frac{\text{sig}(u^{(3)}_i) (1 - \text{sig}(u^{(3)}_i))}{\hat{y}_i (1 - \hat{y}_i)}$$

$$\frac{\partial \mathcal{L}}{\partial u^{(3)}_i} = \frac{1}{N} \left(\frac{(1-y_i)}{(1-\hat{y}_i)} - \frac{y_i}{\hat{y}_i} \right) \cdot \hat{y}_i (1 - \hat{y}_i)$$

$$= \frac{1}{N} \left((1-y_i) \hat{y}_i - y_i (1-\hat{y}_i) \right) = \frac{1}{N} (\hat{y}_i - y_i)$$

$$\begin{aligned} \frac{\partial \text{sig}(u)}{\partial u} &= \frac{\partial \frac{1}{1+e^{-u}}}{\partial u} \\ &= \frac{-1}{(1+e^{-u})^2} \cdot e^{-u} \cdot -1 \\ &= \frac{e^{-u}}{(1+e^{-u})^2} \cdot \frac{1}{(1+e^{-u})} \\ &= \frac{1+e^{-u}-1}{(1+e^{-u})^2} \cdot \frac{1}{(1+e^{-u})} \\ &= \left(1 - \frac{1}{1+e^{-u}}\right) \cdot \frac{1}{1+e^{-u}} \\ &= (1 - \text{sig}(u)) \text{sig}(u) \end{aligned}$$

$$\boxed{\frac{\partial \mathcal{L}}{\partial u^{(3)}_i} = \frac{1}{N} (\hat{y}_i - y_i)} \quad \star$$

$$\frac{\partial \mathcal{L}}{\partial U} \Rightarrow \frac{\partial \mathcal{L}}{\partial U_i} = \frac{\partial \mathcal{L}}{\partial u^{(3)}_j} \cdot \frac{\partial u^{(3)}_j}{\partial U_i}$$

$$u^{(3)}_j = h^{(2)}_{jk} U_k + c$$

solo me importa cuando $k=i$!!

$$\Rightarrow \frac{\partial u^{(3)}_j}{\partial U_i} = \frac{\partial h^{(2)}_{jk} U_k + c}{\partial U_i} = \frac{\partial h^{(2)}_{ji} U_i}{\partial U_i} = h^{(2)}_{ji}$$

$$\frac{\partial \mathcal{L}}{\partial U_i} = \frac{\partial \mathcal{L}}{\partial u^{(3)}_j} \cdot h^{(2)}_{ji} \Rightarrow \boxed{\frac{\partial \mathcal{L}}{\partial U_i} = \frac{\partial \mathcal{L}}{\partial u^{(3)}_j} \cdot h^{(2)}_{ji}} \quad \star$$

$$\frac{\partial \mathcal{L}}{\partial c} = \frac{\partial \mathcal{L}}{\partial u^{(3)}_i} \cdot \frac{\partial u^{(3)}_i}{\partial c} = \frac{\partial \mathcal{L}}{\partial u^{(3)}_i} \cdot 1_i = \text{sum} \left(\frac{\partial \mathcal{L}}{\partial u^{(3)}_i} \right)$$

$$u^{(3)}_i = h^{(2)}_{ik} U_k + c$$

tensor de una dimensión con solo 1_i s

$$\boxed{\frac{\partial \mathcal{L}}{\partial c} = \text{sum} \left(\frac{\partial \mathcal{L}}{\partial u^{(3)}_i} \right)} \quad \star$$

$$\frac{\partial \mathcal{L}}{\partial h^{(2)}} \Rightarrow \frac{\partial \mathcal{L}}{\partial h^{(2)}_{ij}} = \frac{\partial \mathcal{L}}{\partial u^{(3)}_k} \frac{\partial u^{(3)}_k}{\partial h^{(2)}_{ij}} \Rightarrow u^{(3)}_k = h^{(2)}_{k\ell} U_\ell + c$$

solo me da cuando $k=j$

$$\frac{\partial u^{(3)}_k}{\partial h^{(2)}_{ij}} = \frac{\partial h^{(2)}_{k\ell} U_\ell}{\partial h^{(2)}_{ij}} \leftarrow$$

$$= \frac{\partial \mathcal{L}}{\partial u^{(3)}_i} \frac{\partial u^{(3)}_i}{\partial h^{(2)}_{ij}} = \frac{\partial \mathcal{L}}{\partial u^{(3)}_i} \cdot U_j \Rightarrow \frac{\partial \mathcal{L}}{\partial h^{(2)}_{ij}} = \frac{\partial \mathcal{L}}{\partial u^{(3)}_i} U_j$$

$$\Rightarrow \boxed{\frac{\partial \mathcal{L}}{\partial h^{(2)}} = \frac{\partial \mathcal{L}}{\partial u^{(3)}} \otimes U} \quad \star$$

(N, d₂) (N) (d₂) ✓

(← ojo! en Pytorch no hay producto tensorial, pero sí se puede simular con producto de matrices)

$$\frac{\partial \mathcal{L}}{\partial u^{(3)}} \cdot U_{(3, d_2)}$$

(N, 1)

$$\frac{\partial \mathcal{L}}{\partial u^{(2)}} \Rightarrow \frac{\partial \mathcal{L}}{\partial u^{(2)}_{ij}} = \frac{\partial \mathcal{L}}{\partial h^{(2)}_{k\ell}} \frac{\partial h^{(2)}_{k\ell}}{\partial u^{(2)}_{ij}} \leftarrow \text{requiere usando notación de Einstein}$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial u^{(2)}_{ij}} = \frac{\partial \mathcal{L}}{\partial h^{(2)}_{ij}} \cdot \frac{\partial h^{(2)}_{ij}}{\partial u^{(2)}_{ij}}$$

$$h^{(2)}_{k\ell} = \text{sig}(u^{(2)}_{k\ell})$$

$$\Rightarrow \frac{\partial h^{(2)}_{k\ell}}{\partial u^{(2)}_{ij}} \neq 0 \iff k=i \text{ y } j=l$$

$$= \frac{\partial \mathcal{L}}{\partial h^{(2)}_{ij}} \cdot \frac{\partial \text{sig}(u^{(2)}_{ij})}{\partial u^{(2)}_{ij}} = \frac{\partial \mathcal{L}}{\partial h^{(2)}_{ij}} \cdot (\text{sig}(u^{(2)}_{ij}) (1 - \text{sig}(u^{(2)}_{ij})))$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial u^{(2)}} = \frac{\partial \mathcal{L}}{\partial h^{(2)}} \star \text{sig}(u^{(2)}) \star (1 - \text{sig}(u^{(2)}))$$

$$\boxed{\frac{\partial \mathcal{L}}{\partial u^{(2)}} = \frac{\partial \mathcal{L}}{\partial h^{(2)}} \star h^{(2)} \star (1 - h^{(2)})} \quad \star$$

(N, d₂) (N, d₂) (N, d₂) (N, d₂) ✓

ojo! "broadcast" hace la operación con todos los elementos del tensor

$$\frac{\partial \mathcal{L}}{\partial W^{(2)}} \Rightarrow \frac{\partial \mathcal{L}}{\partial W^{(2)}_{ij}} = \frac{\partial \mathcal{L}}{\partial u^{(2)}_{k\ell}} \frac{\partial u^{(2)}_{k\ell}}{\partial W^{(2)}_{ij}} \Rightarrow u^{(2)}_{k\ell} = h^{(1)}_{k\ell} W^{(2)}_{r\ell} + b_\ell$$

$$\Rightarrow \frac{\partial u^{(2)}_{k\ell}}{\partial W^{(2)}_{ij}} = \begin{cases} h^{(1)}_{ki} & \text{si } \ell=j \\ 0 & \text{si } \ell \neq j \end{cases}$$

$$\frac{\partial h^{(1)}_{kr} W^{(2)}_{r\ell} + b_\ell}{\partial W^{(2)}_{ij}} \neq 0 \iff r=i \text{ y } \ell=j$$

en cuyo caso simplificar $h^{(1)}_{ki}$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial W^{(2)}_{ij}} = \frac{\partial \mathcal{L}}{\partial u^{(2)}_{kj}} \cdot h^{(1)}_{ki} \Rightarrow$$

$$\boxed{\frac{\partial \mathcal{L}}{\partial W^{(2)}} = (h^{(1)})^T \cdot \frac{\partial \mathcal{L}}{\partial u^{(2)}}} \quad \star$$

(d₁, d₂) (d₁, N) (N, d₂) ✓

4) Ejemplos Back Prop (continuation)

(5)

$$\frac{\partial \mathcal{L}}{\partial b^{(2)}} \Rightarrow \frac{\partial \mathcal{L}}{\partial b^{(2)}_i} = \frac{\partial \mathcal{L}}{\partial u^{(2)}_{ki}} \cdot \frac{\partial u^{(2)}_{ki}}{\partial b^{(2)}_i}$$

$$u^{(2)}_{kl} = h^{(1)}_{kr} W_{rl}^{(2)} + b^{(2)}_l \leftarrow \text{no se aplica}$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial b^{(2)}_i} = \frac{\partial \mathcal{L}}{\partial u^{(2)}_{ki}} \cdot 1_{ki}$$

$$\frac{\partial \mathcal{L}}{\partial b^{(2)}_i} = \sum_k \frac{\partial \mathcal{L}}{\partial u^{(2)}_{ki}}$$

$$\boxed{\frac{\partial \mathcal{L}}{\partial b^{(2)}} = \text{sum}_1 \left(\frac{\partial \mathcal{L}}{\partial u^{(2)}} \right)} \quad \star$$

como hay varias dimensiones de los espacios, se hace una suma

$$\text{sum}_1 \left(\begin{matrix} \text{matrix} \end{matrix} \right) = \begin{matrix} \text{vector} \end{matrix}$$

$$\frac{\partial \mathcal{L}}{\partial h^{(1)}} \Rightarrow \frac{\partial \mathcal{L}}{\partial h^{(1)}_{ij}} = \frac{\partial \mathcal{L}}{\partial u^{(2)}_{kl}} \cdot \frac{\partial u^{(2)}_{kl}}{\partial h^{(1)}_{ij}}$$

$$u^{(2)}_{kl} = h^{(1)}_{kr} W_{rl}^{(2)} + b^{(2)}_l$$

$$\frac{\partial u^{(2)}_{kl}}{\partial h^{(1)}_{ij}} = \begin{cases} W_{jl}^{(2)} & i=k \\ 0 & i \neq k \end{cases}$$

$$\frac{\partial \mathcal{L}}{\partial h^{(1)}_{ij}} = \frac{\partial \mathcal{L}}{\partial u^{(2)}_{kl}} \cdot W_{jl}^{(2)}$$

$$\boxed{\frac{\partial \mathcal{L}}{\partial h^{(1)}} = \frac{\partial \mathcal{L}}{\partial u^{(2)}} \cdot (W^{(2)})^T} \quad \star$$

$$\frac{\partial \tanh(u)}{\partial u} = \frac{\partial \frac{e^u - e^{-u}}{e^u + e^{-u}}}{\partial u} = 1 - \frac{(e^u - e^{-u})^2}{(e^u + e^{-u})^2} = 1 - (\tanh(u))^2$$

$$\frac{\partial \mathcal{L}}{\partial u^{(1)}} \Rightarrow \frac{\partial \mathcal{L}}{\partial u^{(1)}_{ij}} = \frac{\partial \mathcal{L}}{\partial h^{(1)}_{kl}} \cdot \frac{\partial h^{(1)}_{kl}}{\partial u^{(1)}_{ij}}$$

$$h^{(1)}_{kl} = \tanh(u^{(1)}_{kl})$$

$$\Rightarrow \frac{\partial h^{(1)}_{kl}}{\partial u^{(1)}_{ij}} \neq 0 \Rightarrow k=i, l=j$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial u^{(1)}_{ij}} = \frac{\partial \mathcal{L}}{\partial h^{(1)}_{kl}} \cdot \frac{\partial h^{(1)}_{kl}}{\partial u^{(1)}_{ij}} = \frac{\partial \mathcal{L}}{\partial h^{(1)}_{kl}} \cdot (1 - \tanh(u^{(1)}_{ij})^2)$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial u^{(1)}} = \frac{\partial \mathcal{L}}{\partial h^{(1)}} \star (1 - \tanh(u^{(1)}) \star \tanh(u^{(1)}))$$

$$\Rightarrow \boxed{\frac{\partial \mathcal{L}}{\partial u^{(1)}} = \frac{\partial \mathcal{L}}{\partial h^{(1)}} \star (1 - h^{(1)} \star h^{(1)})} \quad \star$$

$$\boxed{\frac{\partial \mathcal{L}}{\partial W^{(1)}} = \left(\frac{\partial \mathcal{L}}{\partial u^{(1)}} \right)^T \cdot \frac{\partial \mathcal{L}}{\partial u^{(1)}}} \quad \star$$

se obtiene de manera analoga a $\frac{\partial \mathcal{L}}{\partial W^{(2)}}$

$$\boxed{\frac{\partial \mathcal{L}}{\partial b^{(1)}} = \text{sum}_1 \left(\frac{\partial \mathcal{L}}{\partial u^{(1)}} \right)} \quad \star$$

se obtiene de manera analoga a $\frac{\partial \mathcal{L}}{\partial b^{(2)}}$

Recordando:

②

backward

$$\frac{\partial \mathcal{L}}{\partial u^{(3)}} = \frac{1}{N} (\hat{y} - y)$$

$$\frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial u^{(3)}} \cdot h^{(2)}$$

$$\frac{\partial \mathcal{L}}{\partial c} = \text{sum} \left(\frac{\partial \mathcal{L}}{\partial u^{(3)}} \right)$$

$$\frac{\partial \mathcal{L}}{\partial h^{(2)}} = \frac{\partial \mathcal{L}}{\partial u^{(3)}} \otimes U$$

$$\frac{\partial \mathcal{L}}{\partial u^{(2)}} = \frac{\partial \mathcal{L}}{\partial h^{(2)}} * h^{(2)} * (1 - h^{(2)})$$

$$\frac{\partial \mathcal{L}}{\partial W^{(2)}} = (h^{(1)})^T \cdot \frac{\partial \mathcal{L}}{\partial u^{(2)}}$$

$$\frac{\partial \mathcal{L}}{\partial b^{(2)}} = \text{sum}_1 \left(\frac{\partial \mathcal{L}}{\partial u^{(2)}} \right)$$

$$\frac{\partial \mathcal{L}}{\partial h^{(1)}} = \frac{\partial \mathcal{L}}{\partial u^{(2)}} \cdot (W^{(2)})^T$$

$$\frac{\partial \mathcal{L}}{\partial u^{(1)}} = \frac{\partial \mathcal{L}}{\partial h^{(1)}} * (1 - h^{(1)} * h^{(1)})$$

$$\frac{\partial \mathcal{L}}{\partial W^{(1)}} = X^T \cdot \frac{\partial \mathcal{L}}{\partial u^{(1)}}$$

$$\frac{\partial \mathcal{L}}{\partial b^{(1)}} = \text{sum}_1 \left(\frac{\partial \mathcal{L}}{\partial u^{(1)}} \right)$$

Se pueden calcular eficientemente reusando result. todos intermedios y los valores ya calculados en la pasada Forward !!