**Column 1:**

$x_{11} \rightarrow 1$
$x_{21} \rightarrow 2$
$x_{m1} \rightarrow m$

$z_{n \times 1}$

pesos + input + biases $= z_{n \times 1}$

$$W_{n \times m} \, x_{m \times 1} + b_{n \times 1} = z_{n \times 1} \quad (3.1)$$

Usando un batch para entrenar con $p$ ejemplos a la vez:

$$W_{n \times m} \, x_{m \times p} + [\underbrace{b_{n \times 1} \cdots b_{n \times 1}}_{p}] = z_{n \times p}$$

Broadcast $\longrightarrow$

$$W_{n \times m} \, I_{m \times p} + b_{n \times p} = z_{n \times p}$$

batch →    $p$ predicciones a la vez

**Column 2:**

Función de activación softmax a la salida: → probabilidad

$$\hat{y}_k = P(Y = k) = \frac{e^{z_{k1}}}{\sum_j e^{z_{j1}}} \quad (4.1)$$

→ $k = 1, \ldots, n$
→ $j = 1, \ldots, n$

Usando batch:    $(4.2)$

$$\hat{y}_{ki} = P(Y = k \mid X = x_i) = \frac{e^{z_{ki}}}{\sum_j e^{z_{ji}}}$$

→ $k = 1, \ldots, n$
→ $i = 1, \ldots, p$    → $j = 1, \ldots, n$

$$\hat{y} = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

$$\rightarrow \sum_k \hat{y}_{ki} = \frac{\sum_k e^{z_{ki}}}{\sum_j e^{z_{ji}}} = 1 \checkmark$$

$\hat{y}_{ki}$ más alto = predicción del ejemplo $i$.

**Column 3:**

Loss Function para el ejemplo $i$ = Cross entropy = X entropy

$$L_i = -\sum_k y_{ki} \ln(\hat{y}_{ki}) \quad (5.1)$$

↓ valor correcto de la salida

$y_{ki}$ son los elementos de un one-hot vector!

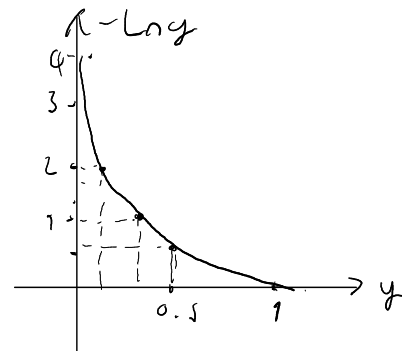$$y_{ki} = 0, 0, \ldots, 0, 1, 0, \ldots, 0$$

↑ $k$ elemnt ↓ clase esperada (etiqueta)

$$y_{ki} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} \\ \\ \\ \leftarrow k \text{ elemnt} \\ \\ \\ \leftarrow n \text{ elemnt} \end{matrix}$$

**Column 4:**

$$\rightarrow L_i = -\ln(\hat{y}_{ki})$$

De la salida esperada para el ej $i$

$$= -\ln\left(\frac{e^{z_{ki}}}{\sum_j e^{z_{ji}}}\right) \quad (5.2)$$

$-\text{Log}$



Cost Function para los pesos, biases del batch:

$$J(W, b) = \frac{1}{p} \sum_i L_i \quad (5.3)$$

$$J(W, b) = \frac{1}{p} \sum_i -\ln\left(\frac{e^{z_{ki}}}{\sum_j e^{z_{ji}}}\right)$$

→ $i = 1, \ldots, p$

Para una sola neurona, una sola entrada:

$x \to \bigcirc \xrightarrow{w} \bigcirc \to z$

$wx + b = z \qquad (6.1)$

$J(w, b)$.

Gráfica computacional: Cada operación es un nodo

$x \to \circledast \to \oplus \to$
$w \nearrow \qquad b \nearrow$

---

Gradient descent:

$$\frac{\partial J}{\partial w} = \lim \frac{J(w+h) - J(w)}{h} \qquad (6.2)$$



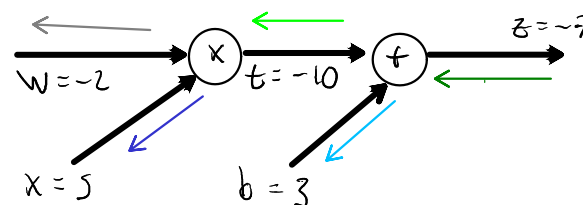$$w = w - \alpha \frac{\partial J}{\partial w} \qquad (6.3)$$

learning rate = step size

De igual manera para $b$:

$$b = b - \alpha \frac{\partial J}{\partial b} \qquad (6.4)$$

---

Backpropagation:

Gráfica computacional:



$z = t + b$
$z = wx + b$

Como es la última (y única) capa de la NN:

$$\frac{dz}{dz} = 1 \qquad (7.1)$$

$$\frac{dz}{db} = \frac{dz}{dz}\frac{dz}{db} = 1$$

$$\frac{dz}{dt} = \frac{dz}{dz}\frac{dz}{dt} = 1$$

$$\frac{dz}{dx} = \frac{dz}{dt}\frac{dt}{dx} = w = -2 \qquad (7.2)$$

---

$$\frac{dz}{dw} = \frac{dz}{dt}\frac{dt}{dw} = x = 5 \qquad (7.3)$$

Ejm: $h = 0.1$
$\uparrow$
cambio de $w$

$w = w_0 + h = -2 + 0.1 = -1.9$

$t = w x_0 = (-1.9)(5) = -9.5$

$z = t + b_0 = (-9.5) + (3) = -6.5$

$z = z_0 + h\frac{dz}{dw}$

$z = -7 + (0.1)5$

$z = -6.5$

**Columna 1**

Ejm: h=0.1
↑ cambio de X

→ $x = x_0 + h = 5 + 0.1 = 5.1$
→ $t = w_0 x = (-2)(5.1) = -10.2$
→ $z = t + b_0 = (-10.2) + (3) = -7.2$

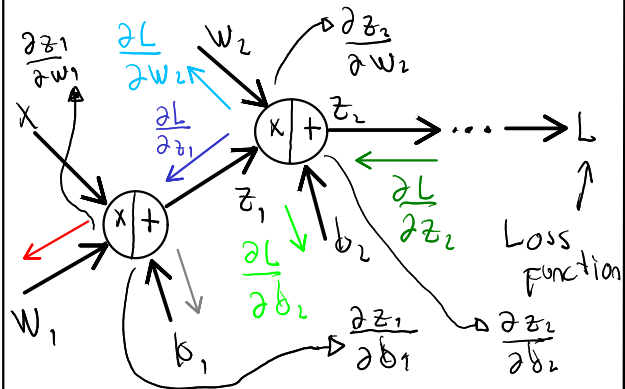→ $z = z_0 + h \frac{dz}{dx}$

$z = -7 + (0.1)(-2)$

$z = -7.2$ ✓

Ejm: h=0.1
↑ cambio de b:

→ $b = b_0 + h = 3 + (0.1) = 3.1$
→ $z = t_0 + b = (-10) + (3.1) = -6.9$

→ $z = z_0 + h \frac{dz}{db}$

$z = -7 + (0.1)(1) = -6.9$ ✓

**Columna 2**

Para dos perceptrones: Gráfica computacional



$\frac{\partial L}{\partial w_2} = ? \quad \frac{\partial L}{\partial b_2} = ? \quad \frac{\partial L}{\partial w_1} = ? \quad \frac{\partial L}{\partial b_1} = ?$

→ $\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2}\frac{\partial z_2}{\partial b_2}$ → $b_2 = b_2 - \alpha \frac{\partial L}{\partial b_2}$

→ $\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2}\frac{\partial z_2}{\partial w_2}$ → $w_2 = w_2 - \alpha \frac{\partial L}{\partial w_2}$

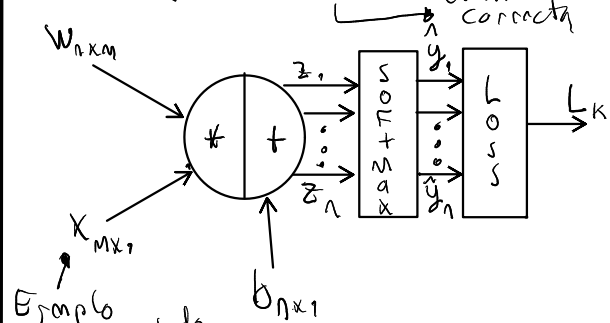→ $\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial z_1}\frac{\partial z_1}{\partial z_2}$

→ $\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial z_1}\frac{\partial z_1}{\partial b_1}$ → $b_1 = b_1 - \alpha \frac{\partial L}{\partial b_1}$

**Columna 3**

→ $\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_1}\frac{\partial z_1}{\partial w_1}$ → $w_1 = w_1 - \alpha \frac{\partial L}{\partial w_1}$

Para un solo ejemplo etiquetado como clase K:

clase correcta



Ejemplo etiquetado

$W_{n\times m}\, X_{m\times 1} + b_{n\times 1} = Z_{n\times 1}$

scores

→ $L = -\sum_K y_K \, Ln\, \hat{y}_K$ ← (5.9)

$(0\ 0\ \dots\ 1\ \dots\ 0)$
K elemento

$L = -Ln\, \hat{y}_K = -Ln\left(\frac{e^{z_K}}{\sum_j e^{z_j}}\right)$

De la salida esperada para el ejemplo

(5.2)

$j = 1,\dots,n$

**Columna 4**

$L = -Ln\left(\frac{e^{z_K}}{\sum_j e^{z_j}}\right) = Ln\sum_j e^{z_j} - z_K$

→ $\frac{\partial L}{\partial z_i} = ?$

$\frac{\partial L}{\partial z_i} = \frac{\partial Ln(\sum_j e^{z_j})}{\partial z_i} - \frac{\partial z_K}{\partial z_i}$

$= \frac{1}{\sum_j e^{z_j}}\frac{\partial}{\partial z_i}\sum_j e^{z_j} - \delta_{iK}$

Delta de kronecker

$= \frac{e^{z_i}}{\sum_j e^{z_j}} - \delta_{iK}$

$\frac{\partial L}{\partial z_i} = \hat{y}_i - y_K$ (10.1)

Predicción para la clase i

Para el score de la salida esperada (clase correcta):

$\frac{\partial L}{\partial z_K} = \hat{y}_K - y_K$ (10.2)

Para inicializar $W$ y $b$ en una NN con pocas capas se puede:

$\rightarrow W_{n \times m} = np.random.randn(n,m) * 0.01$ (11.1)
$\underbrace{\qquad}_{n \times m}$

$\rightarrow b_{n \times 1} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} = np.zeros((n,1))$ (11.2)

$\rightarrow z = W @ x + b$ (11.3)

$\frac{\partial L}{\partial z}$ en forma matricial: $\leftarrow$ (10.2)

$\frac{\partial L}{\partial z_{n \times 1}} = \hat{y}_{n \times 1} - \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}_{n \times 1}$ (11.4) $\leftarrow$ K elemento

$\uparrow$ etiqueta
$\text{one hot}$ $= y_{n \times 1}$ son p6
$\text{vector}$

$\rightarrow \frac{\partial L}{\partial W} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial W} = \frac{\partial L}{\partial z} x$

---

$\frac{\partial L}{\partial W_{n \times m}} = \frac{\partial L}{\partial z_{n \times 1}} \cdot (X_{m \times 1})^T$ (11.5)

$\frac{\partial L}{\partial W_{n \times m}} = (\hat{y}_{n \times 1} - y_{n \times 1}) \cdot (X_{m \times 1})^T$ (11.6)

$= (y\_hat - y) @ X.T$

$\rightarrow \frac{\partial L}{\partial b} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial b} = \frac{\partial L}{\partial z}$ (1)
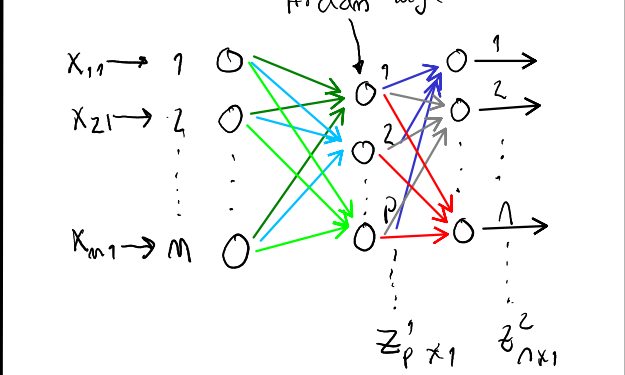
$\frac{\partial L}{\partial b_{n \times 1}} = \frac{\partial L}{\partial z_{n \times 1}} \cdot (1)_{1 \times 1} = \frac{\partial L}{\partial z_{n \times 1}}$

$= \hat{y}_{n \times 1} - y_{n \times 1}$ (11.7)

Después se puede actualizar $W$ y $b$: (6.3)

$\rightarrow W_{n \times m} = W_{n \times m} - \alpha \frac{\partial L}{\partial W_{n \times m}}$ (11.8)
$\searrow$ Ejm: 0.01

$\rightarrow b_{n \times 1} = b_{n \times 1} - \alpha \frac{\partial L}{\partial b_{n \times 1}}$ (11.9)
$\uparrow$ (6.4)

---

Deep learning con dos capas y solo funciones lineales de activación (sin funciones de activación):
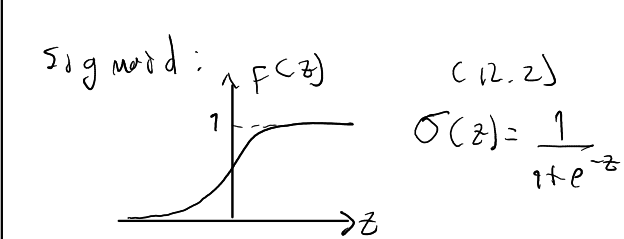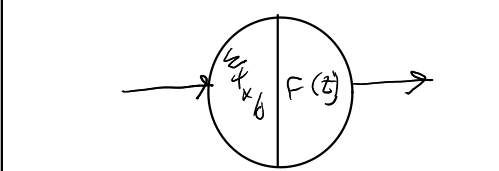

Hidden layer

$z^1_{p \times 1}$   $z^2_{n \times 1}$

$\rightarrow z^1_{p \times 1} = W^1_{p \times m} x_{m \times 1} + b^1_{p \times 1}$

$\rightarrow z^2_{n \times 1} = W^2_{n \times p} z_{p \times 1} + b^2_{n \times 1}$

$= W^2_{n \times p} (W^1_{p \times m} x_{m \times 1} + b^1_{p \times 1}) + b^2_{n \times 1}$

$= (W^2 W^1)_{n \times m} x_{m \times 1} + (W^2 b^1)_{n \times 1} + b^2_{n \times 1}$

$= (W^2 W^1)_{n \times m} x_{m \times 1} + (W^2 b^1 + b^2)_{n \times 1}$ (12.1)

Entonces las hidden layer no actúan, dando igual usar la sola

---

Capa de salida con $n$ perceptrones. Funciones de activación:



sigmoid: (12.2)



$\sigma(z) = \frac{1}{1 + e^{-z}}$

ReLU: Rectified Linear Unit:



$z < 0 \rightarrow F(z) = 0$
$z \geq 0 \rightarrow F(z) = z$

$\max(0, z)$ (12.3)
$\uparrow$
máximo entre $0$ y $z$

Leaky ReLU:



ejm
$z < 0 \rightarrow F(z) = -0.01 z$
$z \geq 0 \rightarrow F(z) = z$
$\max(0.01 z, z)$ (12.4)

**Usando funciones de activación con dos hidden layers:**

R: Ejemplos totales

Usando en Batch de r ejemplos:

$\dfrac{R}{r}$ = pasadas para entrenar con los R elementos

→ Forwardpass
← backpropagation

**\* Forwardpass:**

$$z^1_{p\times r} = W^1_{p\times m} x_{m\times r} + \overbrace{(b^1_{p\times 1} \cdots b^1_{p\times 1})_{p\times r}}^{\text{Broadcast}}$$

Después de aplicar la función de activación de la capa 1:

$$a^1_{p\times r} = F^1(z^1_{p\times r})$$

---

$$\to z^2_{q\times r} = W^2_{q\times p} a^1_{p\times r} + \underbrace{(b^2_{q\times 1} \cdots b^2_{q\times 1})_{q\times r}}_{b^2_{q\times r}}$$

$$\to a^2_{q\times r} = F^2(z^2_{q\times r})$$

$$\to z^3_{n\times r} = W^3_{n\times q} a^2_{q\times r} + b^3_{n\times r} \qquad \text{Broadcast}$$

(11.4)

$\Downarrow$

softmax + loss

etiquetas de los r ejemplos

**\* Backpropagation:**

$$\to \frac{\partial L}{\partial z^3_{n\times r}} = \hat{y}_{n\times r} - y_{n\times r} \quad \leftarrow \text{softmax} \leftarrow (4.2) \quad \text{etiquetas}$$

one hot vector por cada ejemplo:

$$\begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 \cdots 0 \\ \vdots & & \vdots \\ 0 & 0 & 1 \\ 1 & 2 \cdots r \end{pmatrix} \begin{matrix} \leftarrow 1 \\ \leftarrow 2 \\ \\ \leftarrow n \end{matrix}$$

$$\to \frac{\partial L}{\partial a^2} = \frac{\partial L}{\partial z^3} \frac{\partial z^3}{\partial a^2}$$

---

$$\frac{\partial L}{\partial a^2_{q\times r}} = (W^3_{n\times q})^T \frac{\partial L}{\partial z^3_{n\times r}}$$

$$\to \frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial z^3} \frac{\partial z^3}{\partial W^3}$$

$$\frac{\partial L}{\partial W^3_{n\times q}} = \frac{\partial L}{\partial z^3_{n\times r}} (a^2_{q\times r})^T$$

$$\to \frac{\partial L}{\partial b^3} = \frac{\partial L}{\partial z^3} \frac{\partial z^3}{\partial b^3}$$

$$\frac{\partial L}{\partial b^3_{n\times r}} = \frac{\partial L}{\partial z^3_{n\times r}} \cdot 1$$

$$\to \frac{\partial L}{\partial z^2} = \frac{\partial L}{\partial a^2} \frac{\partial a^2}{\partial z^2}$$

$$\frac{\partial L}{\partial z^2_{q\times r}} = \frac{\partial L}{\partial a^2_{q\times r}} \left( \frac{\partial F^2(z^2)}{\partial z^2} \right)$$

$$\to \frac{\partial L}{\partial a^1} = \frac{\partial L}{\partial z^2} \frac{\partial z^2}{\partial a^1}$$

$$\frac{\partial L}{\partial a^1_{p\times r}} = (W^2_{q\times p})^T \frac{\partial L}{\partial z^2_{q\times r}}$$

---

$$\to \frac{\partial L}{\partial W^2} = \frac{\partial L}{\partial z^2} \frac{\partial z^2}{\partial W^2}$$

$$\frac{\partial L}{\partial W^2_{q\times p}} = \frac{\partial L}{\partial z^2_{q\times r}} (a^1_{p\times r})^T$$

$$\to \frac{\partial L}{\partial b^2} = \frac{\partial L}{\partial z^2} \frac{\partial z^2}{\partial b^2}$$

$$\frac{\partial L}{\partial b^2_{q\times r}} = \frac{\partial L}{\partial z^2_{q\times r}} \cdot 1$$

$$\to \frac{\partial L}{\partial z^1} = \frac{\partial L}{\partial a^1} \frac{\partial a^1}{\partial z^1}$$

$$\frac{\partial L}{\partial z^1_{p\times r}} = \frac{\partial L}{\partial a^1_{p\times r}} \left( \frac{\partial F^1(z^1)}{\partial z^1} \right)$$

$$\to \frac{\partial L}{\partial W^1} = \frac{\partial L}{\partial z^1} \frac{\partial z^1}{\partial W^1}$$

$$\frac{\partial L}{\partial W^1_{p\times m}} = \frac{\partial L}{\partial z^1_{p\times r}} (x_{m\times r})^T$$

$$\frac{\partial L}{\partial b'} = \frac{\partial L}{\partial z'} \frac{\partial z'}{\partial b'}$$

$$\frac{\partial L}{\partial b'_{p\times r}} = \frac{\partial L}{\partial z'_{p\times r}} \cdot 1$$

Usando Matriz Jacobiana:

$\frac{\partial L_r}{\partial z^3_n}$ → Loss funtion de r ejemplos

$$= \begin{pmatrix} \frac{\partial L_1}{\partial z^3_1} & \cdots & \frac{\partial L_1}{\partial z_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial L_r}{\partial z^3_1} & \cdots & \frac{\partial L_r}{\partial z^3_n} \end{pmatrix} = \frac{\partial L}{\partial z_{r\times n}}$$

$$\frac{\partial L}{\partial z^3_{r\times n}} = \hat{y}_{r\times n} - y_{r\times n}$$

one hot vectors!

$$\begin{pmatrix} 0 & 1 & \cdots & 0 \\ 1 & 0 & \cdots & 0 \\ \vdots & & \ddots & \\ 0 & 0 & \cdots & 1 \end{pmatrix}$$

---

$$\frac{\partial L_r}{\partial a^2_q} = \frac{\partial L_r}{\partial z^3_n} \frac{\partial z^3_n}{\partial a^2_q}$$

$$\frac{\partial L}{\partial a^2_{r\times q}} = \frac{\partial L}{\partial z^3_{r\times n}} \frac{\partial z^3}{\partial a^2_{n\times q}}$$

$$= (\hat{y}_{r\times n} - y_{r\times n}) W^3_{n\times q}$$

$$\frac{\partial L_r}{\partial W^3_{n\times q}} = \frac{\partial L_r}{\partial z^3_n} \frac{\partial z^3_n}{\partial W^3_{n\times q}}$$

$$\frac{\partial L}{\partial W^3_{r\times n\times q}} = (\hat{y}_{r\times n} - y_{r\times n}) a^2_q$$

→ Producto tensorial?

$$\frac{\partial L_r}{\partial b^3_n} = \frac{\partial L_r}{\partial z^3_n} \frac{\partial z^3_n}{\partial b^3_n}$$

$$\frac{\partial L}{\partial b^3_{r\times n}} = \frac{\partial L}{\partial z^3_{r\times n}} I_{n\times n} = \frac{\partial L}{\partial z^3_{r\times n}}$$

$$\frac{\partial L_r}{\partial W^2_{q\times p}} = \frac{\partial L_r}{\partial z^3_n} \frac{\partial z^3_n}{\partial a^2_q} \frac{\partial a^2_q}{\partial z^2_q} \frac{\partial z^2_q}{\partial W^2_{q\times p}}$$

---

$$= (\hat{y}_{r\times n} - y_{r\times n}) W^3_{n\times q} \left(\frac{\partial F^2(z^2)}{\partial z^2}\right)_{q\times q} a^2_p$$

$\underbrace{\qquad}_{r\times q}$

$$= \frac{\partial L}{\partial W^2_{r\times q\times p}} \qquad \rightarrow \frac{\partial L_r}{\partial z^2_q}$$

$$\frac{\partial L_r}{\partial b^2_q} = \frac{\partial L_r}{\partial z^3_n} \frac{\partial z^3_n}{\partial a^2_q} \frac{\partial a^2_q}{\partial z^2_q} \frac{\partial z^2_q}{\partial b^2_q} \quad \rightarrow I_{q\times q}$$
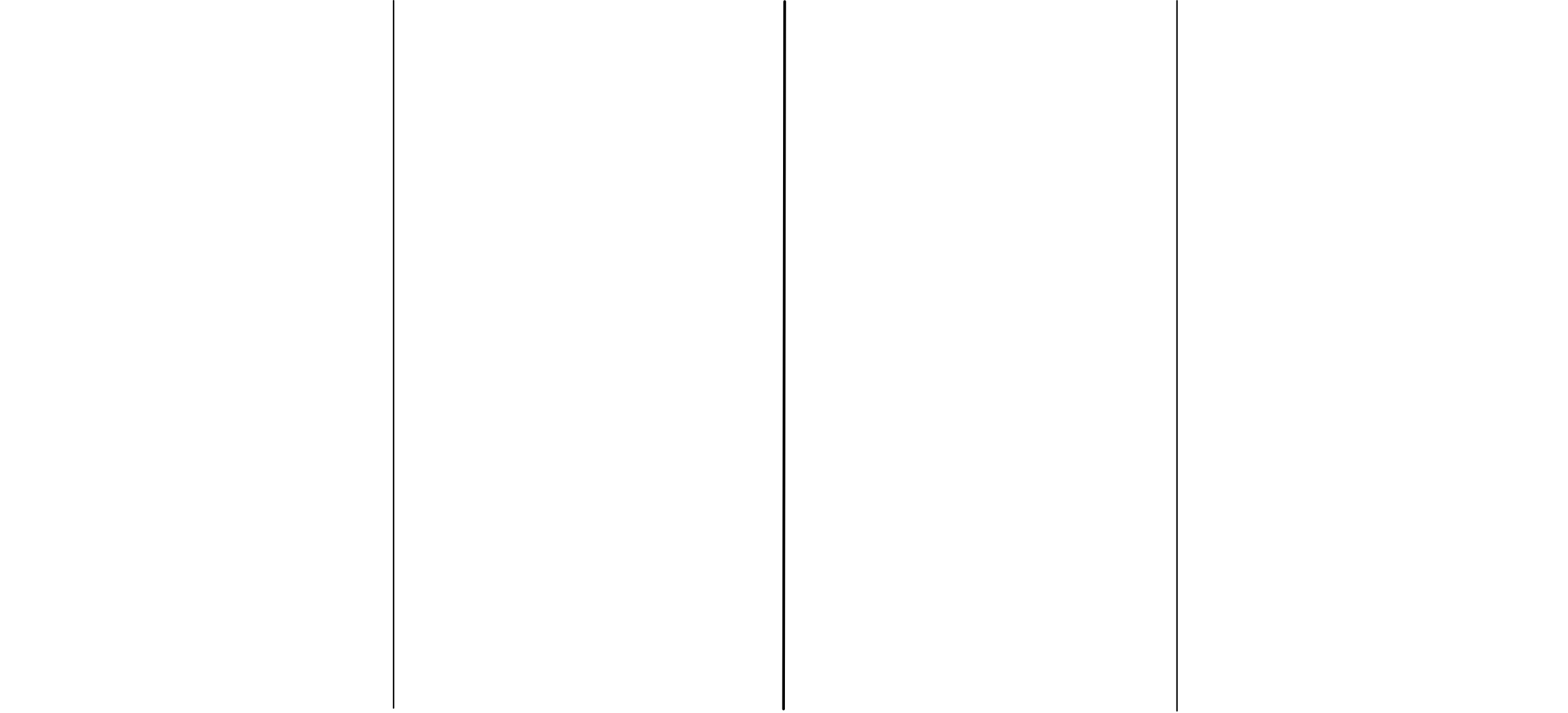
$$= \frac{\partial L}{\partial b^2_{r\times q}}$$

$$\frac{\partial L_r}{\partial W^1_{p\times m}} = \frac{\partial L_r}{\partial z^3_n} \frac{\partial z^3_n}{\partial a^2_q} \frac{\partial a^2_q}{\partial z^2_q} \frac{\partial z^2_q}{\partial a^1_p}$$

$$\cdot \frac{\partial a^1_p}{\partial z^1_p} \frac{\partial z^1_p}{\partial W^1_{p\times m}}$$

$$\frac{\partial L}{\partial W^1_{r\times p\times m}} = \underbrace{\frac{\partial L_r}{\partial z^2_q} W^2_{q\times p} \left(\frac{\partial F^1(z^1)}{\partial z^1}\right)_{p\times p}}_{\frac{\partial L_r}{\partial z^1_p}} x_m$$
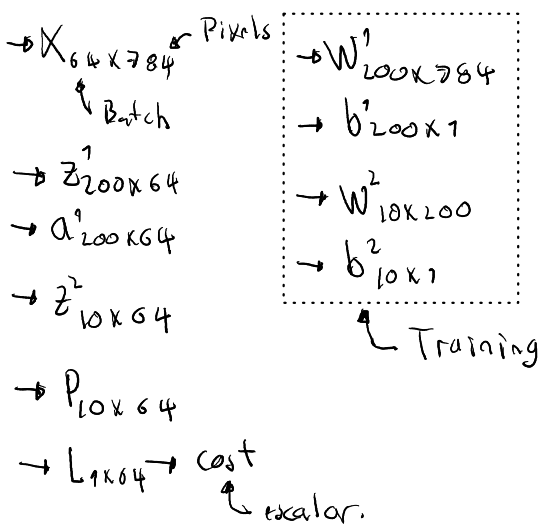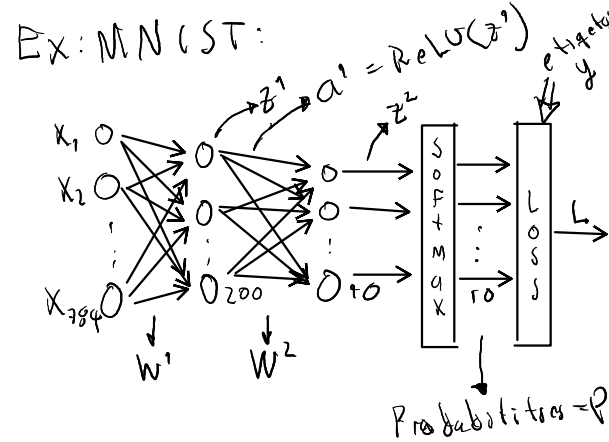
---

$$\frac{\partial L_r}{\partial b'_p} = \frac{\partial L_r}{\partial z'_p} \frac{\partial z'_p}{\partial b'_p}$$

$\rightarrow I_{p\times p}$

$$= \frac{\partial L}{\partial b'_{r\times p}}$$

# EX: MNIST:



$X_1, X_2, \ldots, X_{784}$ inputs → $z^1$ → $a^1 = ReLU(z^1)$ → $z^2$ → SOFTMAX → LOSS → $L$

etiqueta $y$

Probabilities = $P$

$W^1$, $W^2$

---

→ $X_{64 \times 784}$ — Pixels, Batch

→ $z^1_{200 \times 64}$

→ $a^1_{200 \times 64}$

→ $z^2_{10 \times 64}$

→ $P_{10 \times 64}$

→ $L_{1 \times 64}$ → cost → escalar.

Training box:
→ $W^1_{200 \times 784}$
→ $b^1_{200 \times 1}$
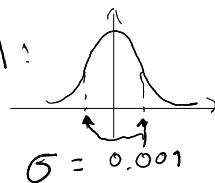→ $W^2_{10 \times 200}$
→ $b^2_{10 \times 1}$

---

## Inicialmente!

$$W^1_{200 \times 784} = np.random.randn(200,784) \cdot 0.001$$

Distribución normal:

$$\sigma = 0.001$$

→ $b^1_{200 \times 1} = np.zeros((200,1))$

→ $W^2_{10 \times 200} = np.random.randn(10,200) \cdot 0.001$

→ $b^2_{10 \times 1} = np.zeros(((10,1))$

## Forward:

$$z^1_{200 \times 64} = W^1_{200 \times 784}\left(X_{64 \times 784}\right)^T + \left(b^1\, b^1 \cdots b^1\right)_{200 \times 64}$$

(1 … 64)

$$a^1_{200 \times 64} = ReLU\left(z^1_{200 \times 64}\right)$$

$$z^2_{10 \times 64} = W^2_{10 \times 200}\, a^1_{200 \times 64} + \left(b^2\, b^2 \cdots b^2\right)_{10 \times 64}$$

---

$$P_{10 \times 64} = \frac{\left(e^{z^2}\right)_{10 \times 64}}{\left(\sum_j e^{z^2_{ji}}\right)_{1 \times 64}} = \hat{y}_{10 \times 64} \quad (4.1)$$

suma cada columna (ej $i$) de $\left(e^{z^2}\right)_{10 \times 64}$

$(5.2)$

$$\left(e^L\right)_{1 \times 64} = \left(P_{a_1},\, P_{b_2} \cdots P_{z_{64}}\right)$$

etiqueta ej. 1    etiqueta ej 2

→ $L_{1 \times 64} = Ln\left(e^L\right)_{1 \times 64}$

→ $cost = \dfrac{\sum\limits_i^{64} -L_{1i}}{64}$

$(5.3)$

---

## Backward:

→ $\dfrac{\partial L}{\partial z^2}_{10 \times 64} = P_{10 \times 64} - y_{10 \times 64}$

$\hat{y}_{10 \times 64}$

etiquetas

one hot vector por cada ejemplo:

$$\begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ & \cdots & \\ 0 & 0 & 1 \end{pmatrix} \begin{matrix} \leftarrow 1 \\ \leftarrow 2 \\ \\ \leftarrow 10 \end{matrix}$$

$1\ 2 \cdots 64$

→ $\dfrac{\partial L}{\partial W^2} = \dfrac{\partial L}{\partial z^2}\dfrac{\partial z^2}{\partial W^2}$

$$\dfrac{\partial L}{\partial W^2_{10 \times 200}} = \dfrac{\dfrac{\partial L}{\partial z^2_{10 \times 64}}\left(a^1_{200 \times 64}\right)^T}{64}$$

→ $\dfrac{\partial L}{\partial b^2} = \dfrac{\partial L}{\partial z^2}\dfrac{\partial z^2}{\partial b^2}$  $\color{red}{1}$

$$\dfrac{\partial L}{\partial b^2_{10 \times 1}} = \begin{pmatrix} \sum\limits_j^{64} \dfrac{\partial L}{\partial z^2_{1j}} \\ \\ \sum\limits_j^{64} \dfrac{\partial L}{\partial z^2_{10j}} \end{pmatrix}_{10 \times 1} \Big/ 64$$

suma clases de cada fila

$$\rightarrow \frac{\partial L}{\partial a^1} = \frac{\partial L}{\partial z^2} \frac{\partial z^2}{\partial a_1}$$

$$\frac{\partial L}{\partial a^1_{200 \times 64}} = \left( W^2_{10 \times 200} \right)^T \frac{\partial L}{\partial z^2_{10 \times 64}}$$

$$\rightarrow \frac{\partial L}{\partial z^1} = \frac{\partial L}{\partial a^1} \frac{\partial a^1}{\partial z^1}$$

$$\frac{\partial L}{\partial z^1_{200 \times 64}} = \frac{\partial L}{\partial a^1_{200 \times 64}}$$