

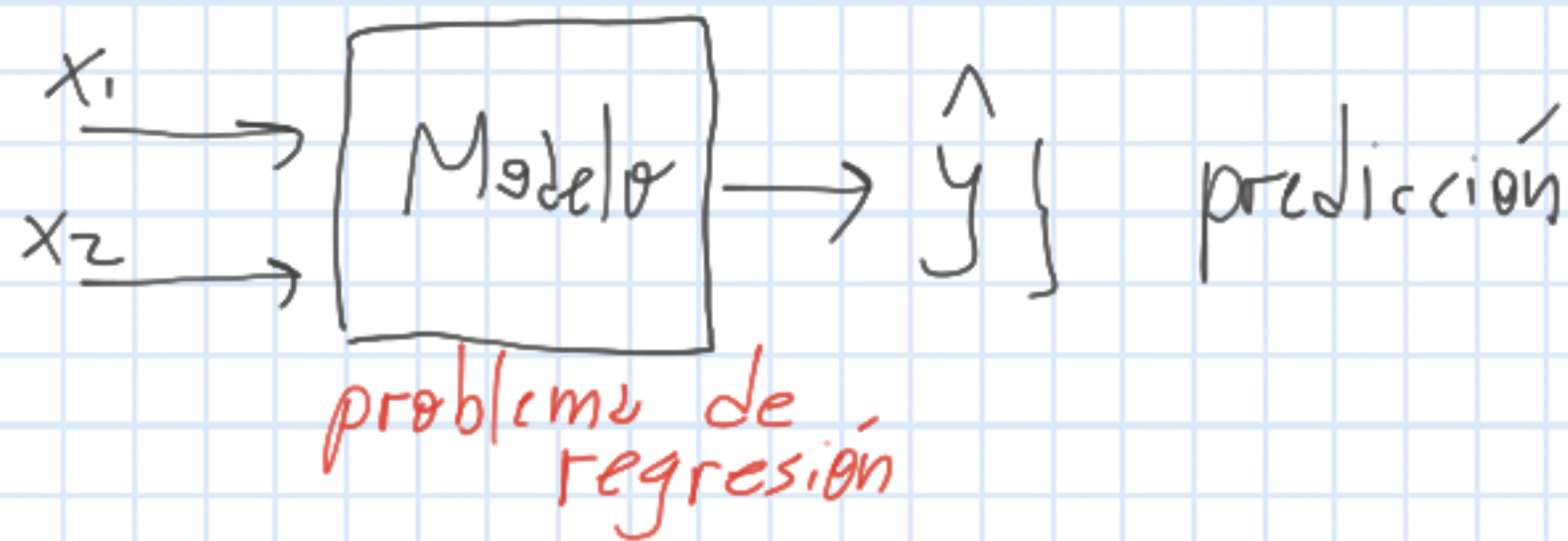
¿Por qué necesitamos modelos no lineales?

$n=4$

x_1	x_2	$y = x_1 \text{ xor } x_2$
0	0	0
0	1	1
1	0	1
1	1	0

$m=2$

label



* Modelo: $\hat{y}_i = x_{i,1} w_1 + x_{i,2} w_2 + b$ \rightarrow ¿parámetro desconocida? $\rightarrow w_1, w_2, b$

escalar

* Loss function $J(\bar{w}) = \frac{1}{4} \sum_{i=1}^4 (y_i - \hat{y}_i)^2$; $\bar{w} = [w_1, w_2, b]^T$

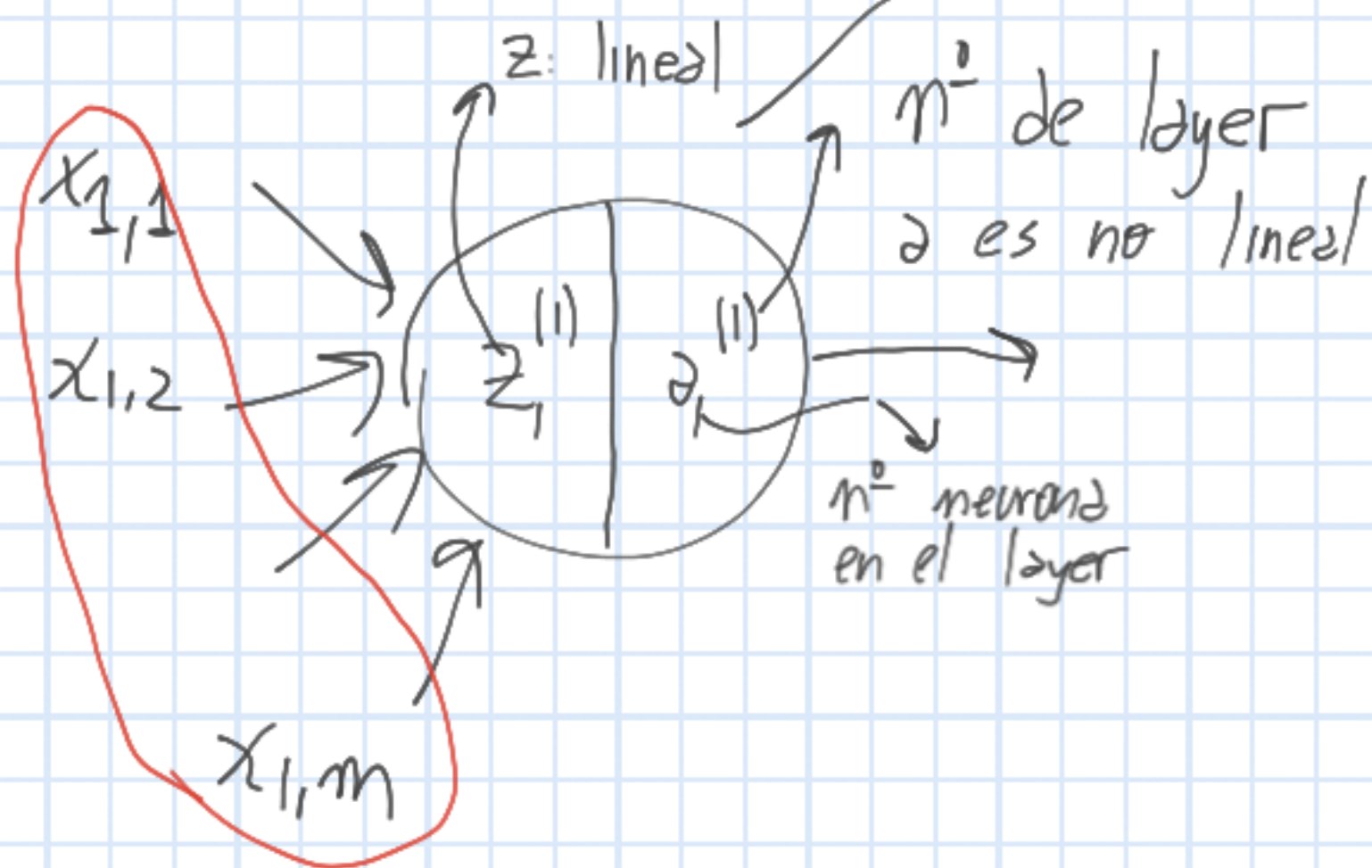
* Optimizador: $\nabla_{\bar{w}} J = 0 \rightarrow \begin{bmatrix} \partial J / \partial w_1 \\ \partial J / \partial w_2 \\ \partial J / \partial b \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \bar{w} = \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix} = (X^T X)^{-1} X^T Y$

$\begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix}$

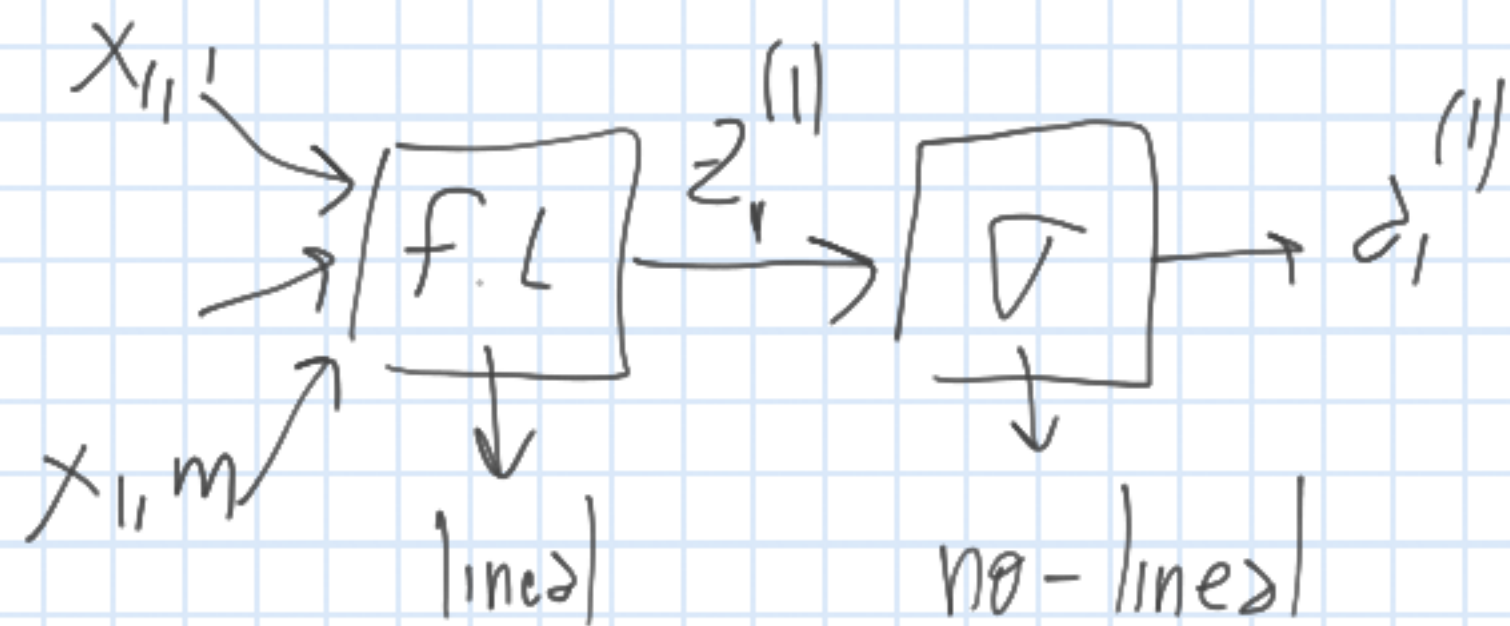
$\begin{bmatrix} 0 \\ 0 \\ 1/2 \end{bmatrix}$

Redes Neuronales

* Neurona



→ La fila 1 de mi data set

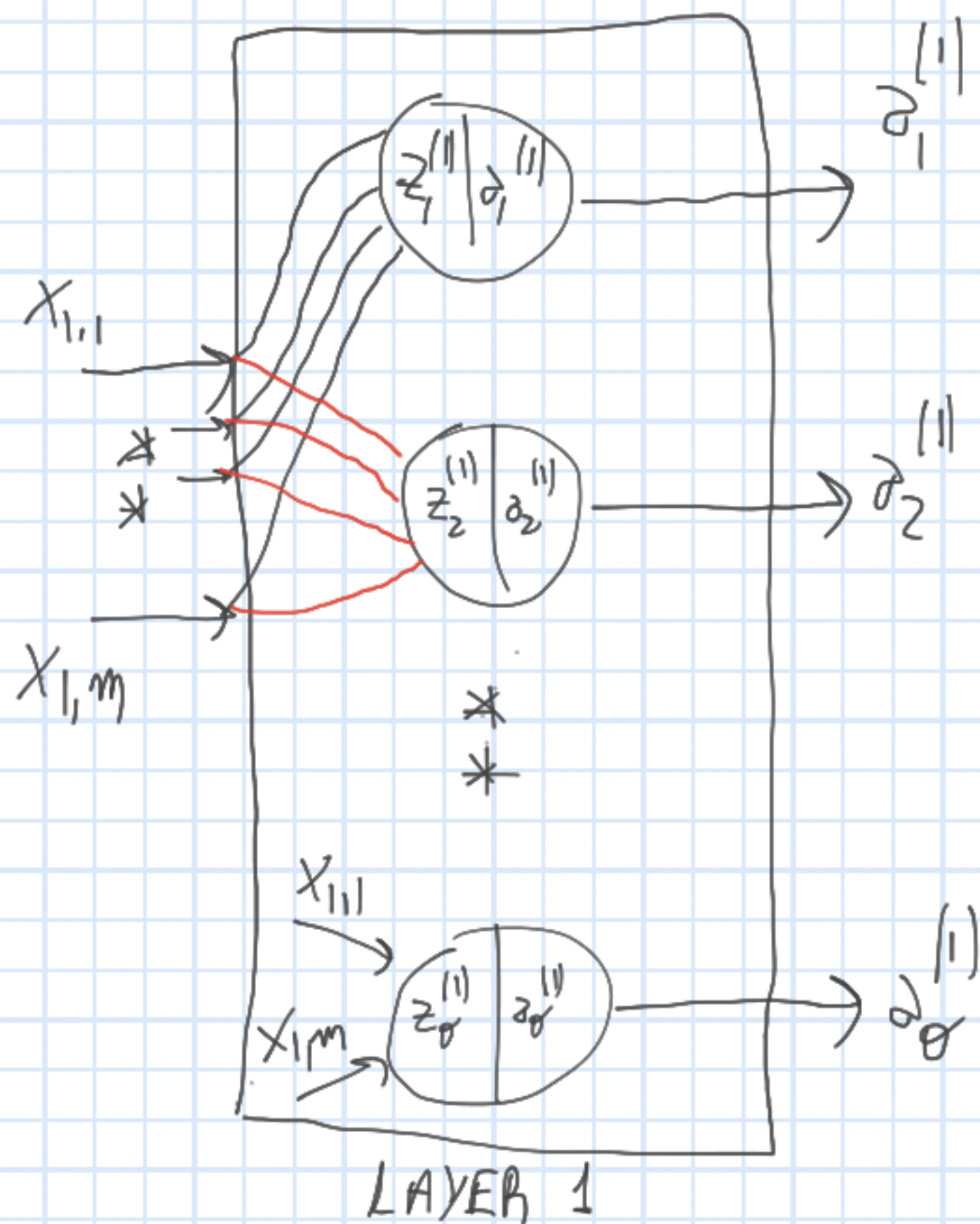


$$z_1^{(1)} = w_{1,1}^{(1)} x_{1,1} + \dots + w_{1,m}^{(1)} x_{1,m} + b_1^{(1)}$$

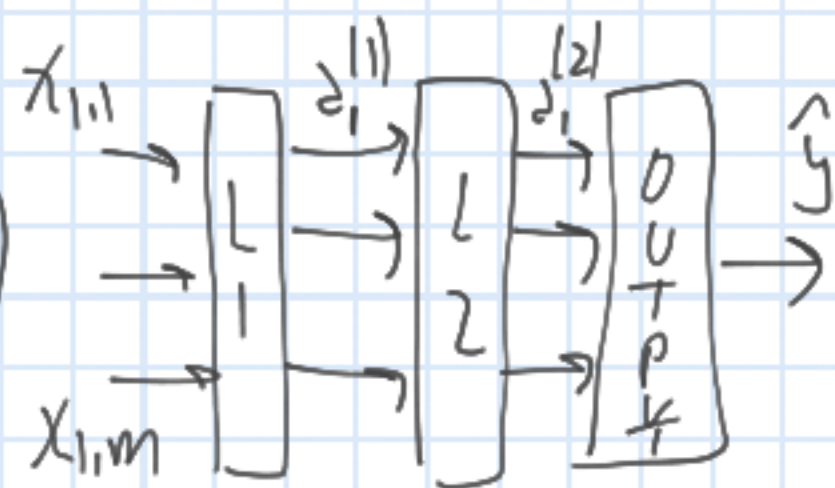
$$a_1^{(1)} = \sqrt{z_1^{(1)}} \rightarrow \text{"activación"}$$

$\left. \begin{array}{l} \text{sigmoid} \\ \text{ReLU} \\ \text{tanh} \end{array} \right\} \text{SOFTMAX}$

* Definición de layer



Feedforward



DATASET

$$\left\{ \begin{array}{l} x_{1,1} \dots x_{1,m} \rightarrow \bar{x}_1 \\ \vdots \\ x_{n,1} \dots x_{n,m} \end{array} \right\} \text{features}$$

$$a_1^{(1)} = \sigma(z_1^{(1)}) = \sigma(x_{1,1}w_{1,1} + \dots + x_{1,m}w_{1,m} + b_1^{(1)})$$

$$a_{\theta}^{(1)} = \sigma(z_{\theta}^{(1)}) = \sigma(x_{1,1}w_{\theta,1} + \dots + x_{1,m}w_{\theta,m} + b_{\theta}^{(1)})$$

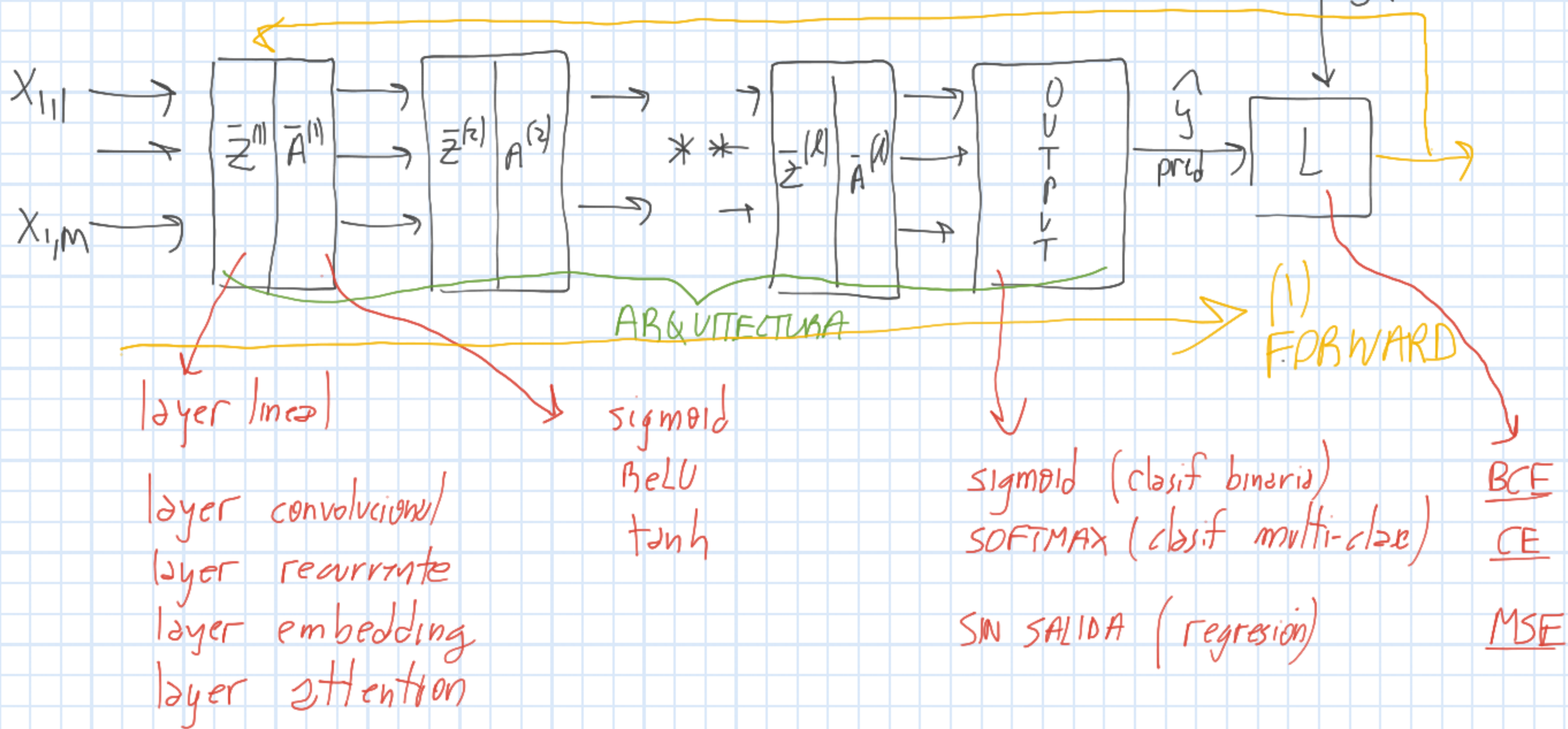
$$\bar{A}^{(1)} = \sigma(\bar{Z}^{(1)})$$

$$\bar{Z}^{(1)} = \underbrace{W^{(1)}}_{\theta \times m} \underbrace{\bar{X}_1}_{m \times 1} + \underbrace{b^{(1)}}_{\theta \times 1}$$

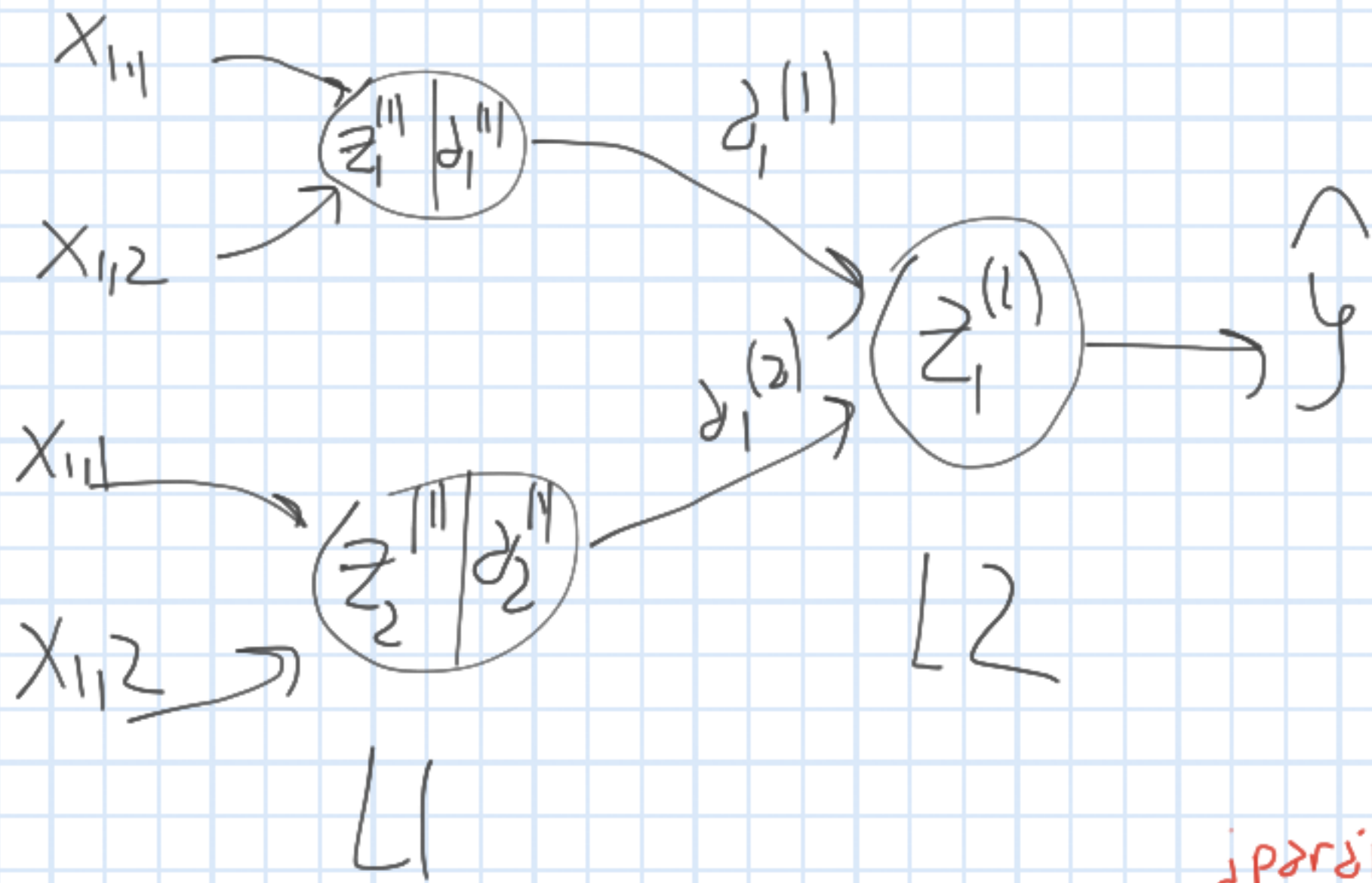
¿parameters?

$$(\theta \times m) + \theta \times 1 = \theta \times (m+1)$$

(2) Backpropagation



* Resolver XOR con redes nev.



¿parámetros?
 $3+3+3=9$

* Forward

$$z_{i,1}^{(1)} = x_{i,1} w_{1,1}^{(1)} + x_{i,2} w_{1,2}^{(1)} + b_1^{(1)}$$

$$a_{i,1}^{(1)} = \sigma(z_{i,1}^{(1)}) = \frac{1}{1 + e^{-z_{i,1}^{(1)}}}$$

$$z_{i,2}^{(1)} = x_{i,1} w_{2,1}^{(1)} + x_{i,2} w_{2,2}^{(1)} + b_2^{(1)}$$

$$a_{i,2}^{(1)} = \sigma(z_{i,2}^{(1)})$$

$$z_{i,1}^{(2)} = a_{i,1}^{(1)} w_{1,1}^{(2)} + a_{i,2}^{(1)} w_{2,1}^{(2)} + b_1^{(2)}$$

$$\hat{y}_i = z_{i,1}^{(2)}$$

$$L_i = (y_i - \hat{y}_i)$$

* Backpropagation

$$\frac{\partial L}{\partial w_{1,1}^{(1)}} =$$

*
*

$$\frac{\partial L}{\partial w_{2,1}^{(1)}} =$$

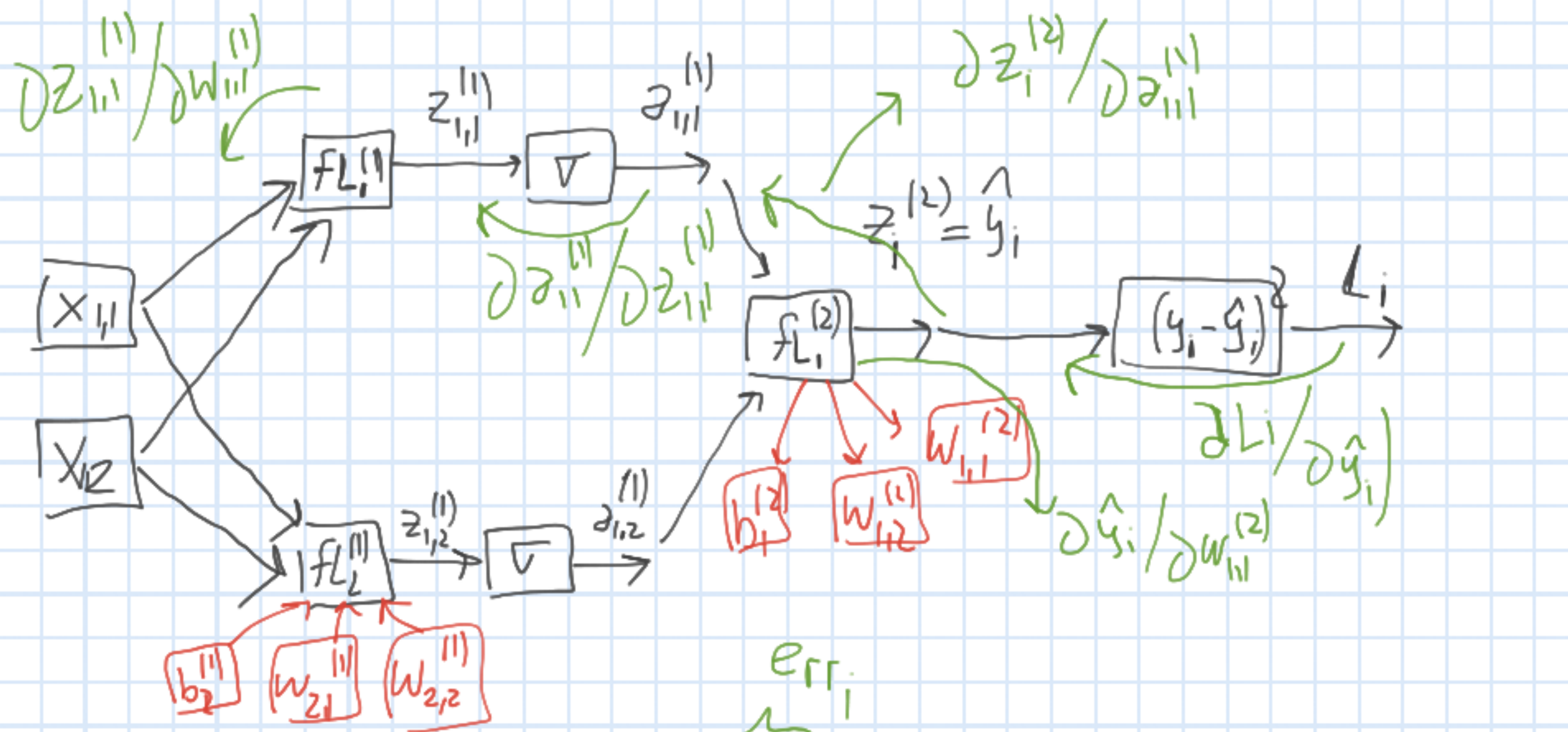
$$\frac{\partial L}{\partial b_1^{(1)}} =$$

$$\frac{\partial L}{\partial b_2^{(1)}} =$$

*
*

$$\frac{\partial L}{\partial w_{1,2}^{(1)}} =$$

9 equations



$$\frac{\partial L}{\partial w_{1,1}^{(2)}} = \underbrace{\left(\frac{\partial L}{\partial \hat{y}_1} \right)}_{err_i} \underbrace{\left(\frac{\partial \hat{y}_1}{\partial w_{1,1}^{(2)}} \right)}_z = 2 (y_1 - \hat{y}_1) (-1) \frac{z_{1,1}^{(2)}}{z}$$

$$\frac{\partial L}{\partial w_{1,1}^{(1)}} = \frac{\partial L}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial z_{1,1}^{(2)}} \frac{\partial z_{1,1}^{(2)}}{\partial z_{1,1}^{(1)}} \frac{\partial z_{1,1}^{(2)}}{\partial w_{1,1}^{(1)}} = -2 err_i w_{1,1}^{(2)} \sigma(z_{1,1}^{(1)}) (1 - \sigma(z_{1,1}^{(1)})) * x_{1,1}$$

$$\bar{W} = \begin{pmatrix} w_{11}^{(1)} \\ w_{12}^{(1)} \\ w_{21}^{(1)} \\ w_{22}^{(1)} \\ b_1^{(1)} \\ b_2^{(1)} \\ w_{11}^{(2)} \\ w_{12}^{(2)} \\ b_1^{(2)} \end{pmatrix} = \text{Inicializaci3n en valores random}$$

Iteraci3n usando SGD

$$\bar{W}_{t+1} \leftarrow \bar{W}_t - \alpha \nabla_{\bar{W}} J$$

learning rate α calculamos $\nabla_{\bar{W}} J$

$$\begin{pmatrix} w_{11}^{(1)} \\ \vdots \\ b_1^{(2)} \end{pmatrix}_{t+1} \leftarrow \begin{pmatrix} w_{11}^{(1)} - \alpha \partial L / \partial w_{11}^{(1)} \\ \vdots \\ b_1^{(2)} - \alpha \partial L / \partial b_1^{(2)} \end{pmatrix}_t$$

Algoritmo SGD

* Inicializar η params random

* For epoch in range(100):

for i in range(4):

- (1) Forward $X_1 = X[i,1], X_2 = X[i,2], y_i = Y[i] \rightarrow \hat{y}_i$
- (2) $err = y_i - \hat{y}_i$
- (3) Backpropagation $\rightarrow \eta$ der. parciales
- (4) Update pesos $\rightarrow \bar{W} = \bar{W} - \alpha D(3)$

(5) Calcular $MSE = \frac{1}{4} \sum (y - \hat{y})^2$

$w_{11}^{(1)}, w_{12}^{(1)}, \dots$

itero
por cada muestra

