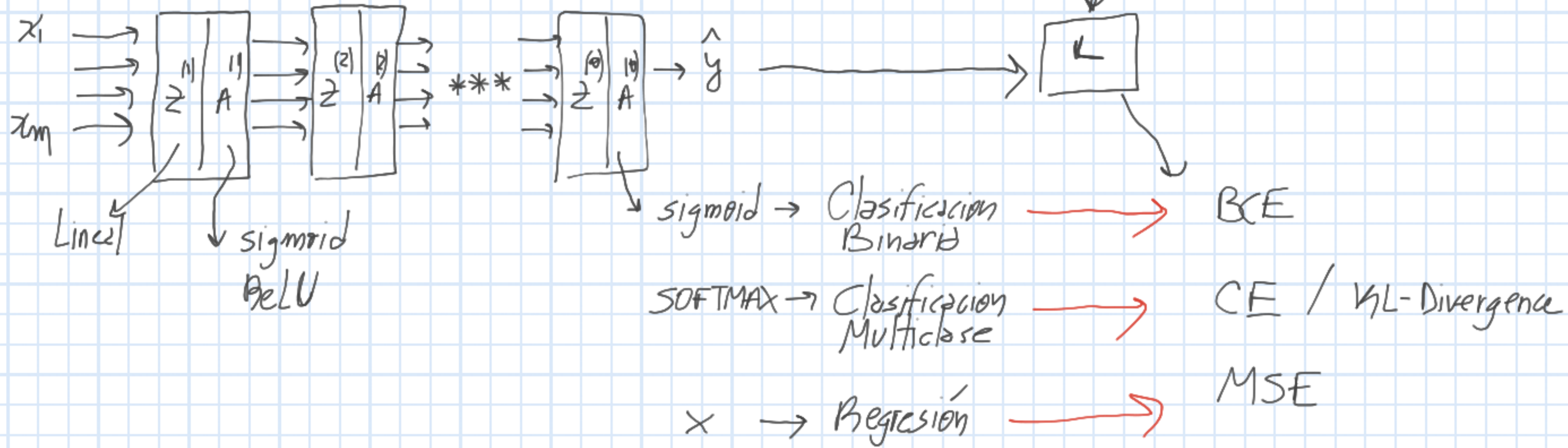
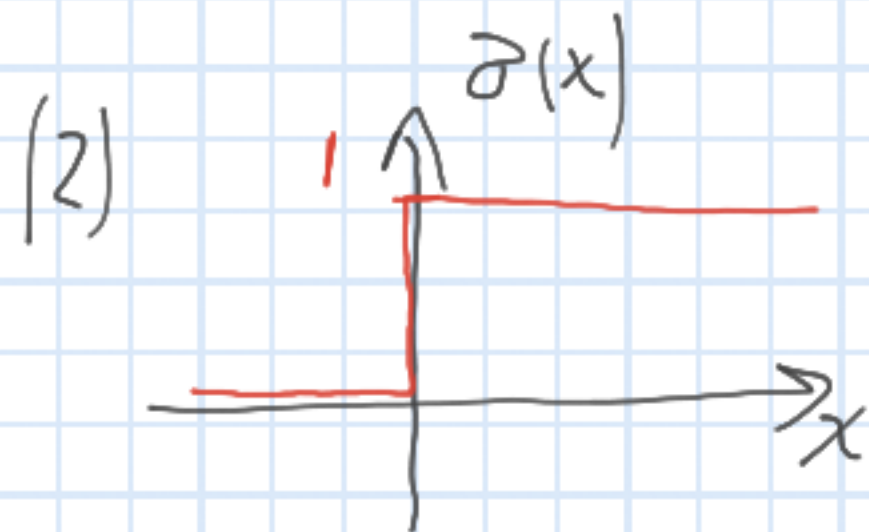


Funciones de activación



(1) No use función de activación

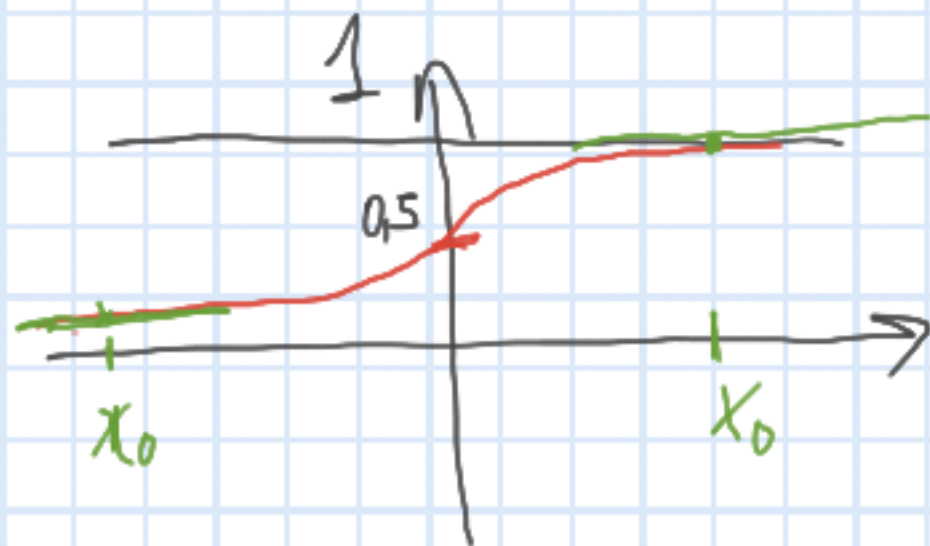
$$x \rightarrow \boxed{z^{(1)}} \rightarrow \boxed{z^{(2)}} \rightarrow \hat{y} \rightarrow \hat{y} = f(x) = \underbrace{z^{(2)}}_{\text{Lin}} \left(\underbrace{z^{(1)}(x)}_{\text{Lin}} \right)$$



$$\sigma(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases} \rightarrow \sigma'(x) = \begin{cases} 0 & x > 0 \\ 0 & x \leq 0 \end{cases}$$

$$\textcircled{w_i} \leftarrow w_i - \alpha \frac{\partial L}{\partial w_i}$$

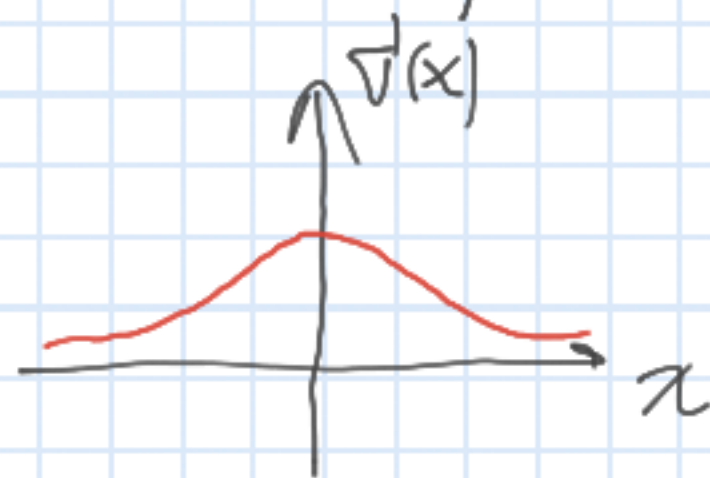
(3) $\sigma(x) = \Gamma(x)$ sigmoid



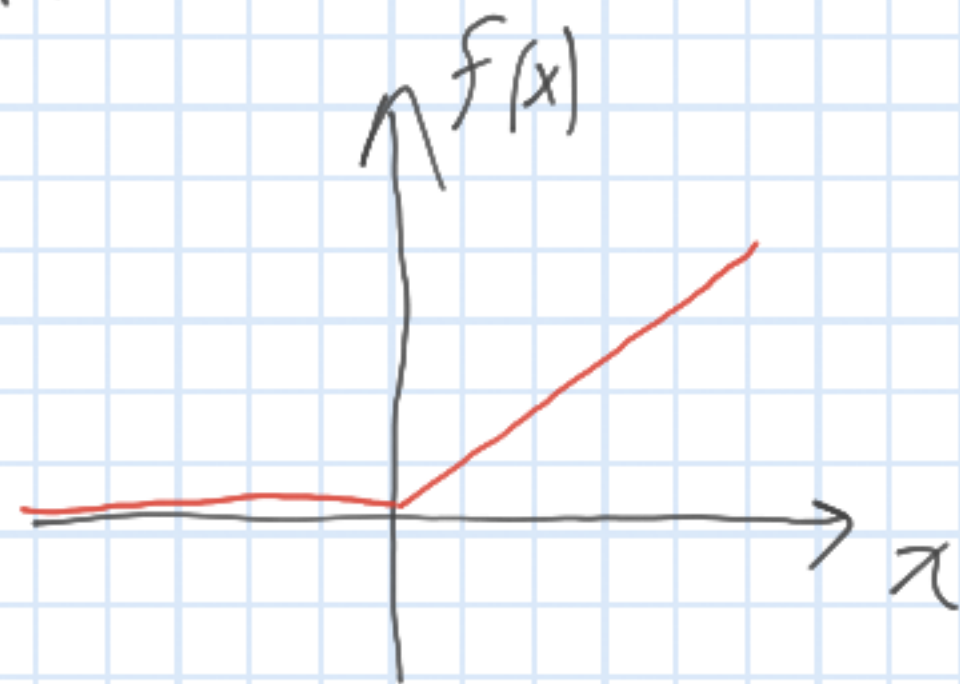
$$\Gamma(x) = \frac{1}{1 + e^{-x}}$$

$$\Gamma'(x) = \Gamma(x)(1 - \Gamma(x))$$

¿Problema? \rightarrow Optimización lenta
 \rightarrow Vanishing Gradients



(4) ReLU

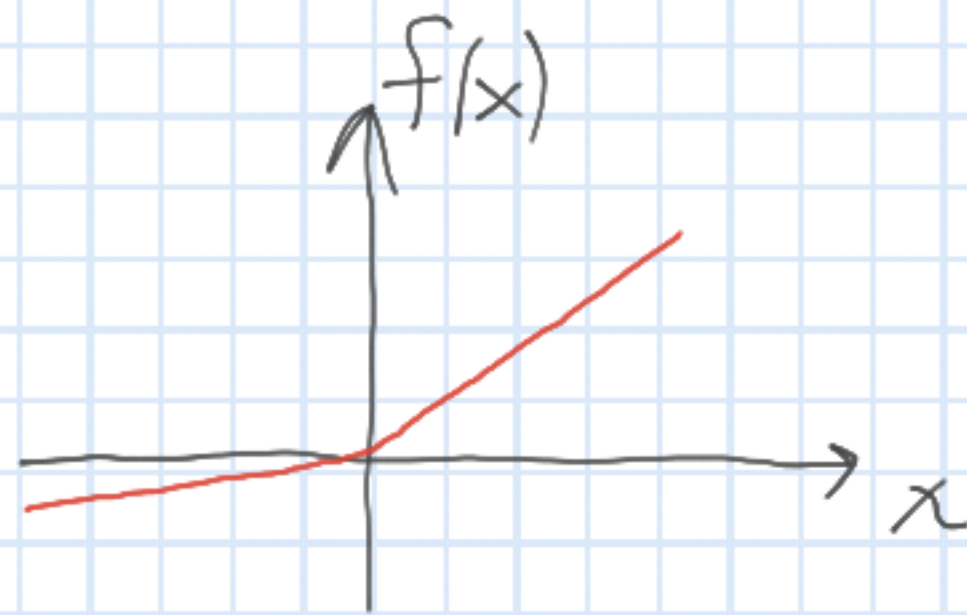


$$f(x) = \begin{cases} x & x > 0 \\ 0 & x \leq 0 \end{cases} = \max(x, 0)$$

$$f'(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$$

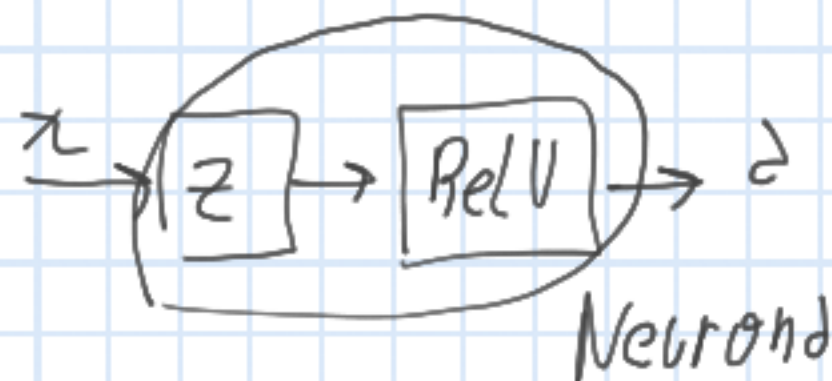
Para valores negativos de x la derivada es cero y mata a la neurona

(5) Leaky ReLU

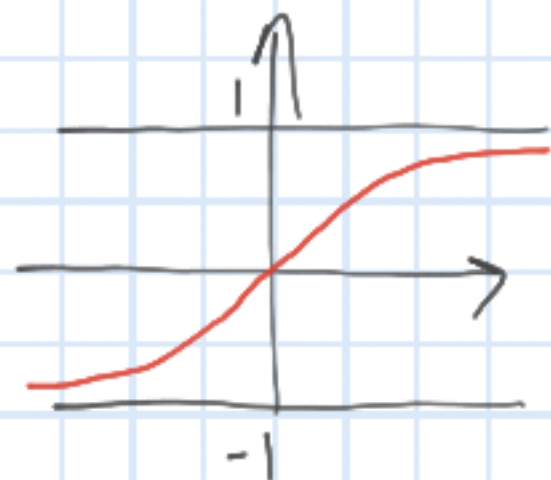


$$f(x) = \begin{cases} x & x > 0 \\ 0.01x & x \leq 0 \end{cases}$$

$$f'(x) = \begin{cases} 1 & x > 0 \\ 0.01 & x \leq 0 \end{cases}$$



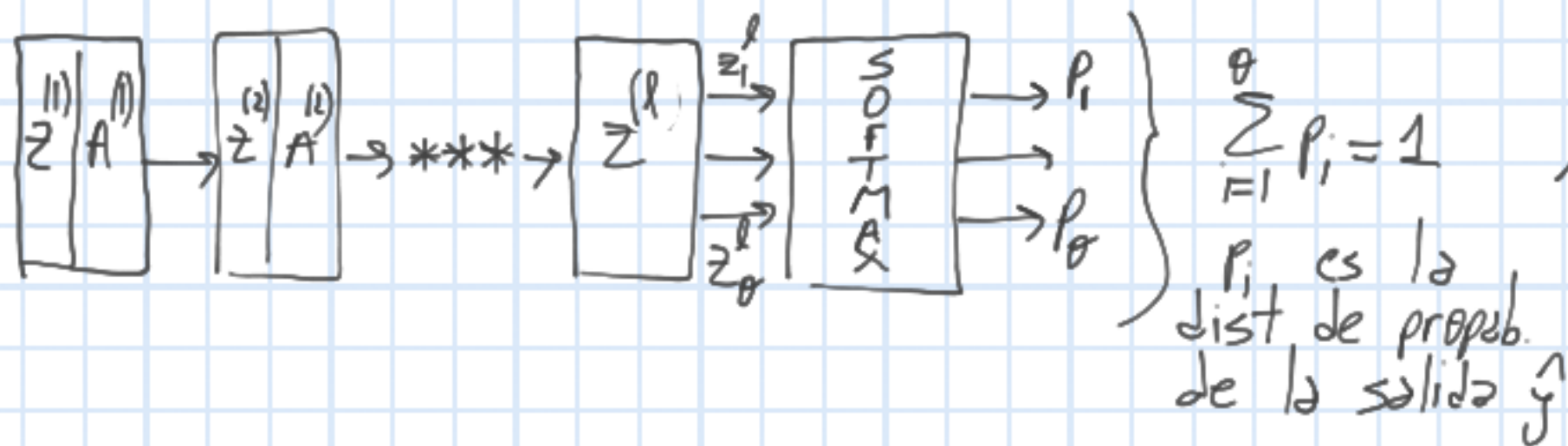
$$(6) \tanh_{(x)} = 2\sigma(x) - 1$$



Funciones de salida

- Regresión → ϕ (ninguna) ✓
- Clasificación binaria → σ (sigmoid) ✓
- Clasificación multiclase → SOFTMAX

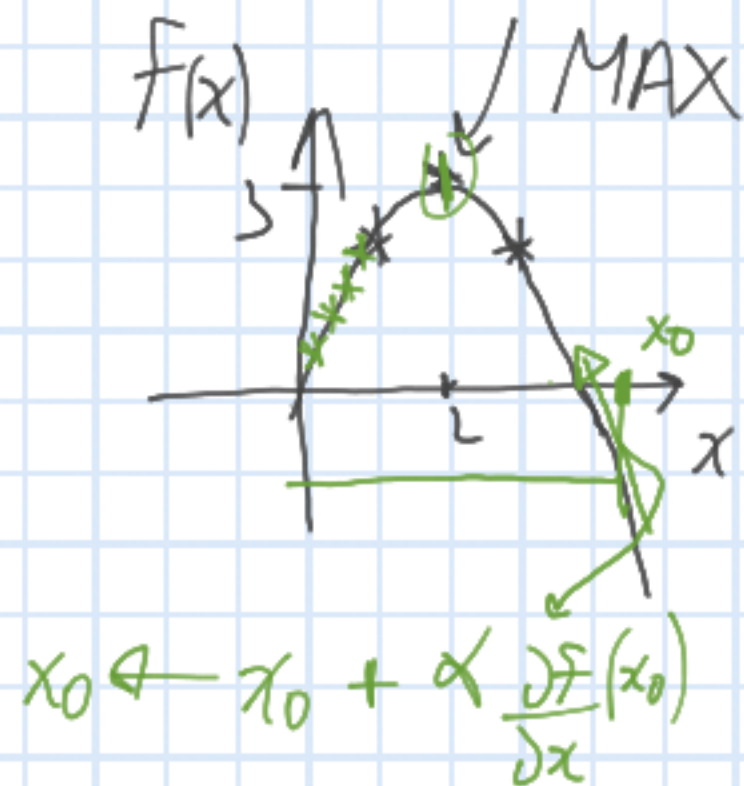
Clasificador de θ clases



$$p_1 = \frac{e^{z_1}}{\sum_{i=1}^{\theta} e^{z_i}}$$

$$p_{\theta} = \frac{e^{z_{\theta}}}{\sum_{i=1}^{\theta} e^{z_i}}$$

Optimización en DL



$$f(x) = -(x-2)^2 + 3$$

Implementar
GD para encontrar
el máximo de $f(x)$

$$\hat{y} = \text{NNET}(x)$$

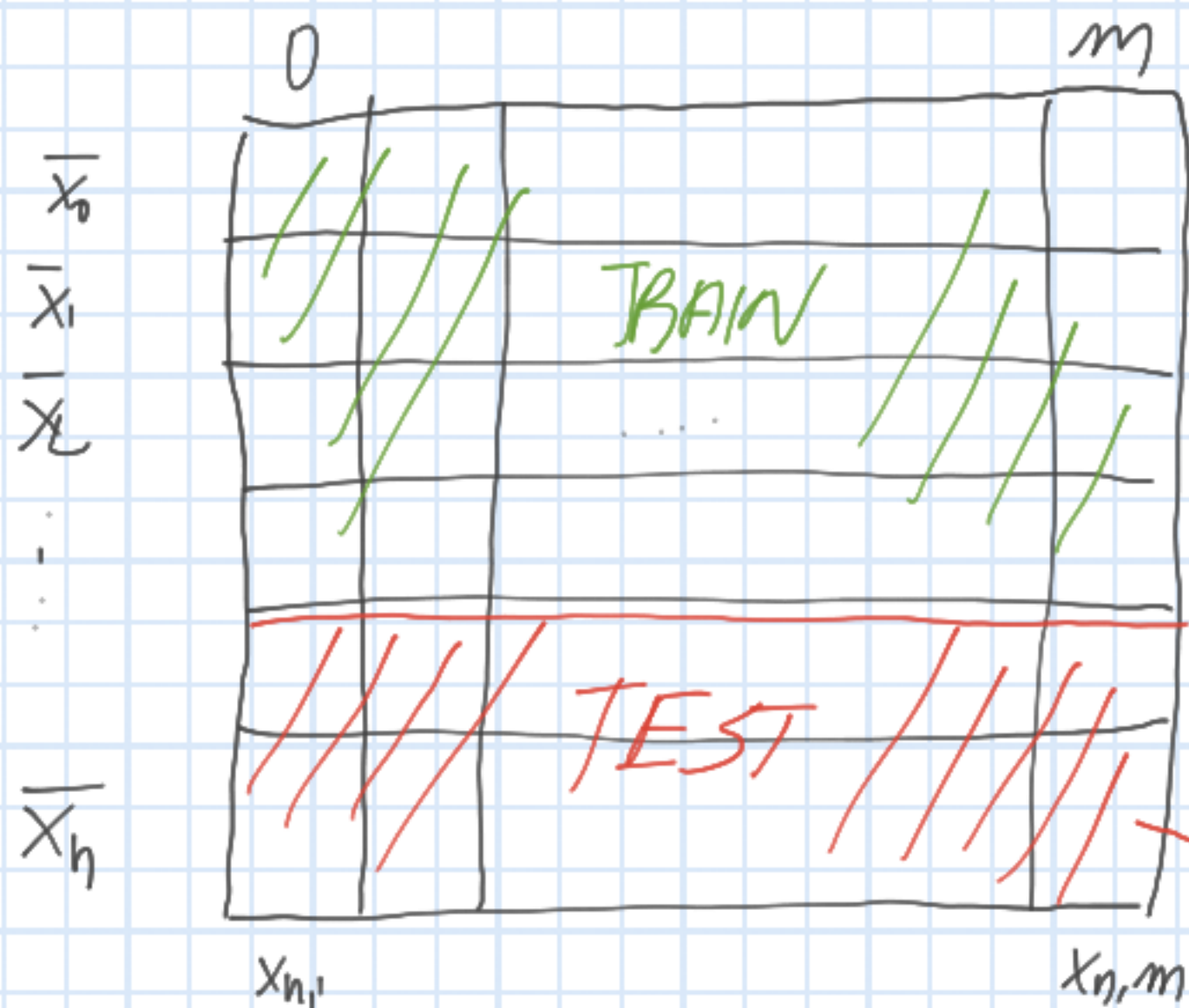
$$L = f(\bar{x}, \hat{y}, y)$$

$\hat{y} = \bar{g}(\bar{w})$ params

Optimización en DL

- ¿Cómo se trabaja con el dataset? (1)
 - ¿Tengo que normalizar los datos? (2)
 - Vanishing gradients (3)
 - Exploding gradients (4)
 - Algoritmos de optimización (5)
 - Regularización (6)
 - L2
 - L1
 - Dropout
- } Evitar overfitting

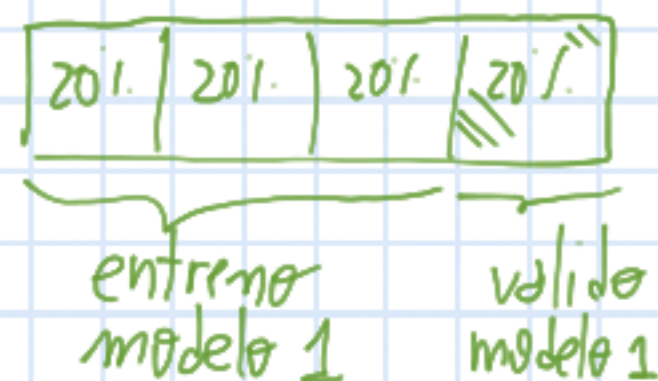
(1) Dataset en deep learning



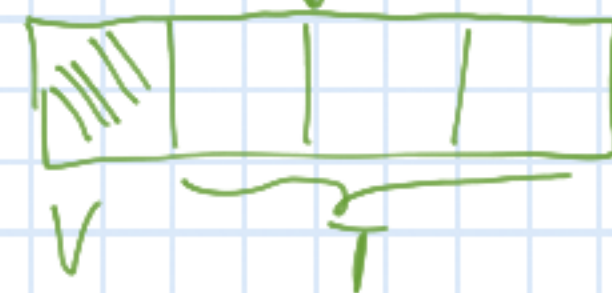
DL \rightarrow m y n son muy grandes

ML tradicional

Cross Validation \rightarrow HP search



AUC modelo 1



AUC modelo 4

AUC promedio

En deep learning

Calcular métricas performance de mi modelo

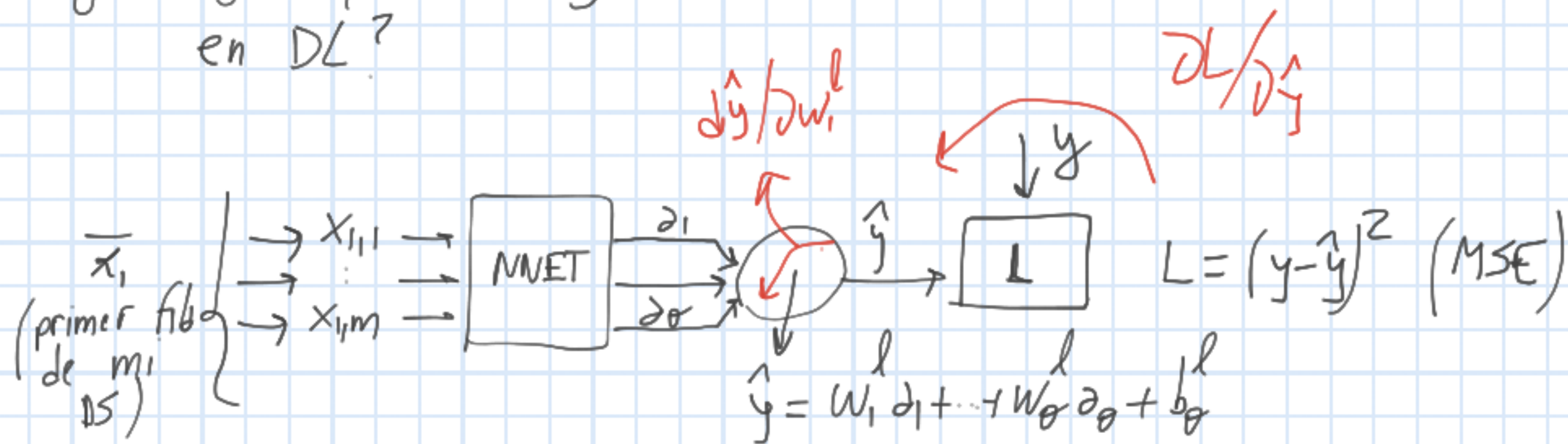


Clasificación:

Reg \rightarrow RMSE

F1
P
R
ACC
AUC

(Z) ¿Tengo que normalizar datos en DL?



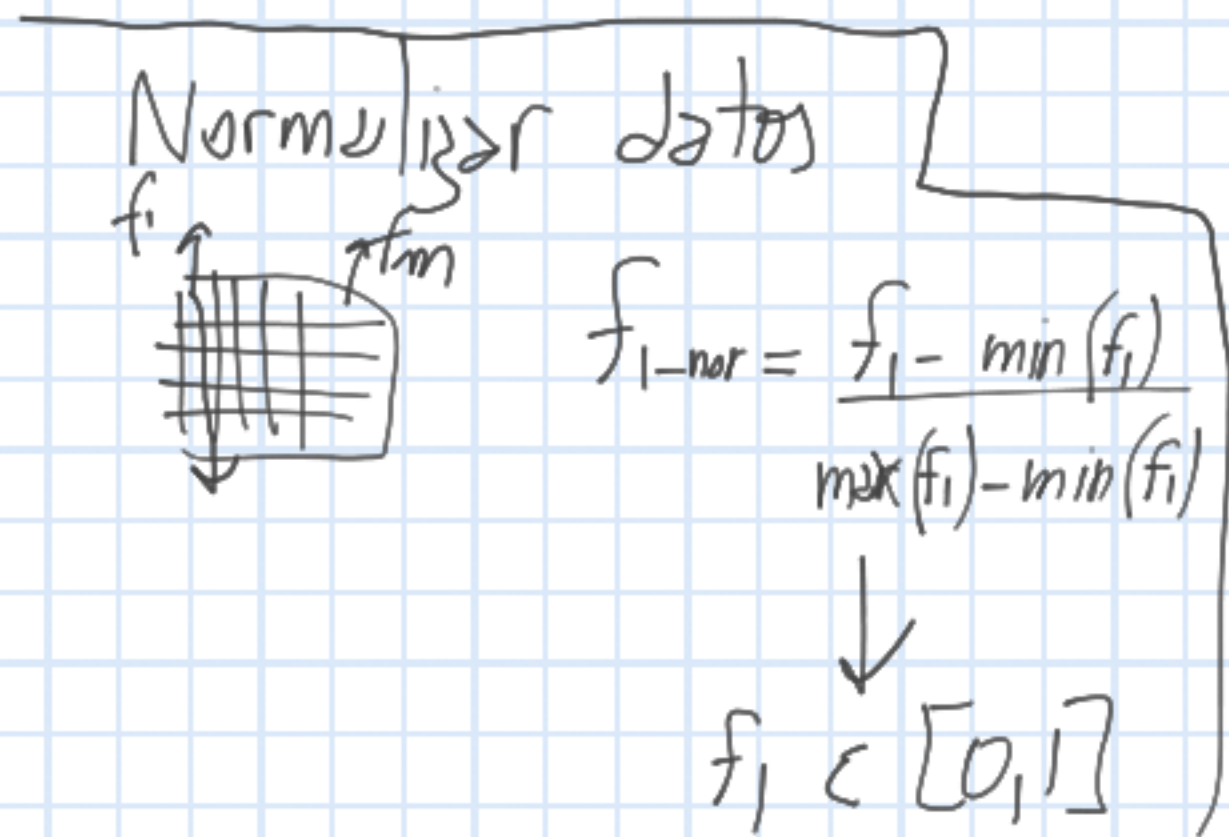
$$\frac{\partial L}{\partial w_1^l} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_1^l} = -2(y - \hat{y}) a_1$$

$$w_1^l \leftarrow w_1^l - \alpha (-2(y - \hat{y})) a_1$$

$$w_1^l \leftarrow w_1^l + \alpha 2(y - \hat{y}) a_1$$

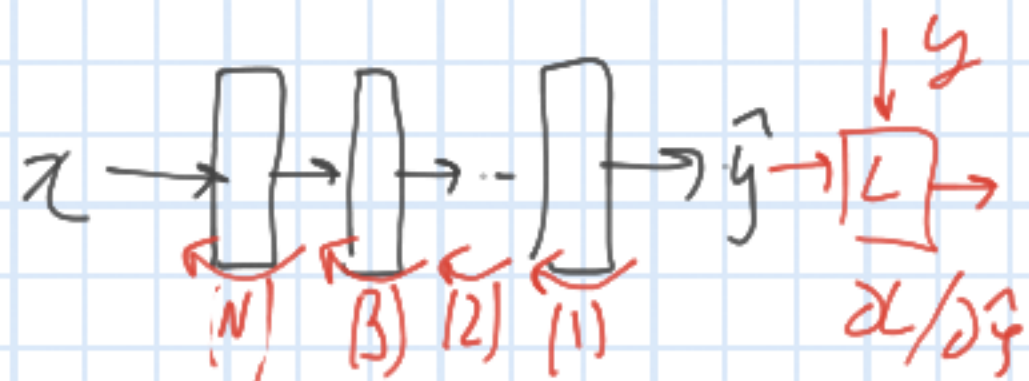
$$w_2^l \leftarrow w_2^l + \alpha 2(y - \hat{y}) a_2$$

Sol \rightarrow normalizar datos



(3) y (4)

Vanishing Gradients

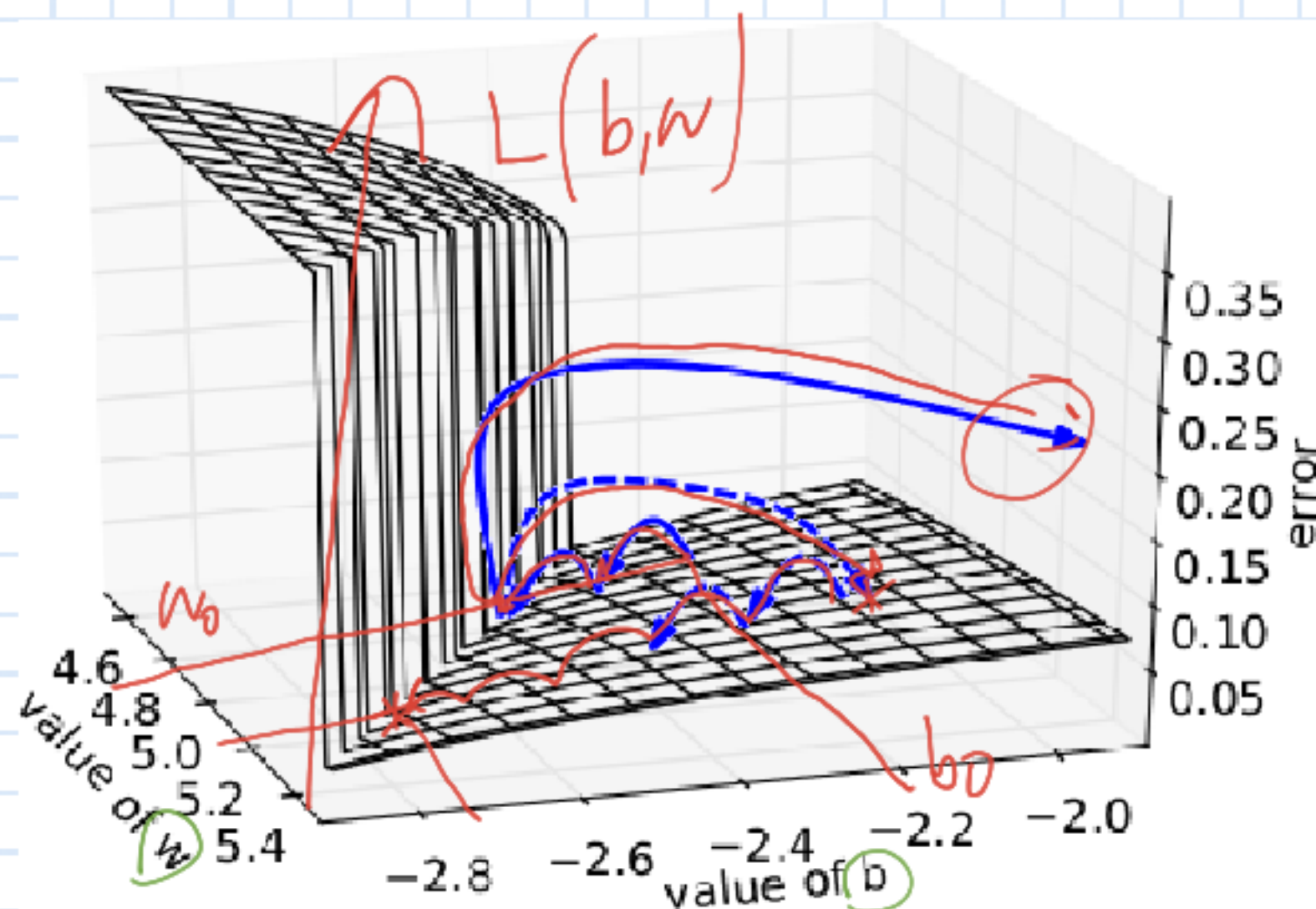


$$\frac{\partial \hat{y}}{\partial w_j^{(1)}} = \frac{\partial L}{\partial \hat{y}} (1) (2) (3) \dots (n)$$

Prede ser cerea
si tenga muchos
capas

ReLU

Exploding Gradients

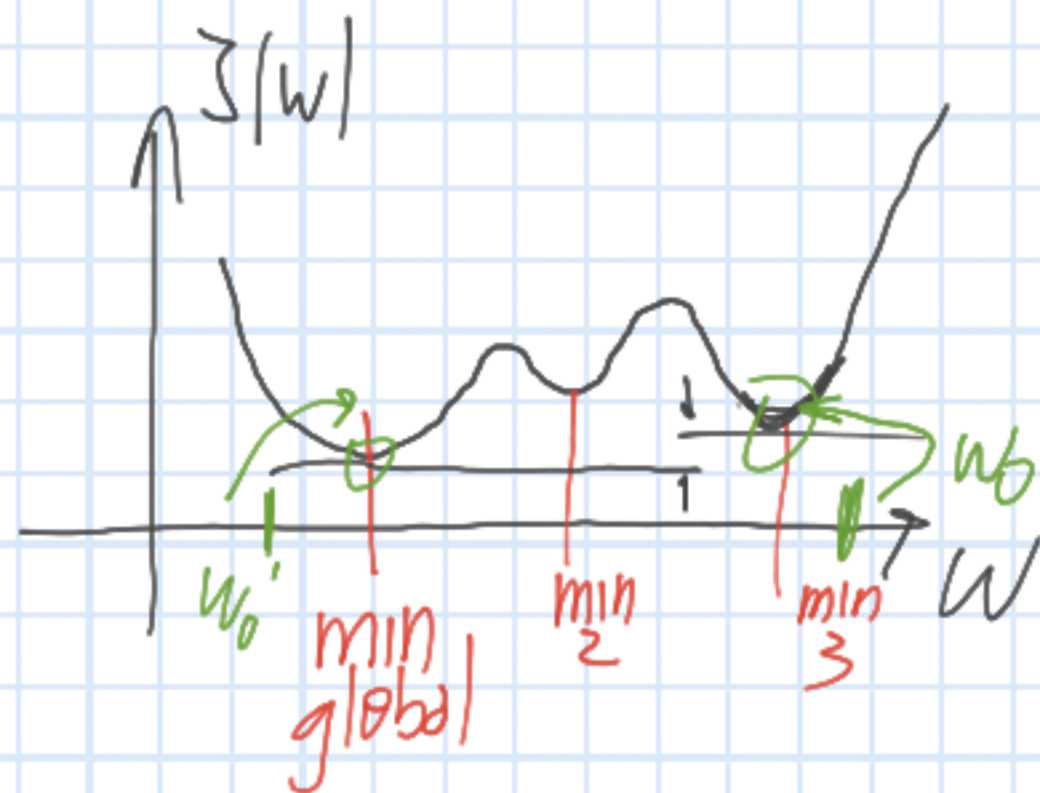
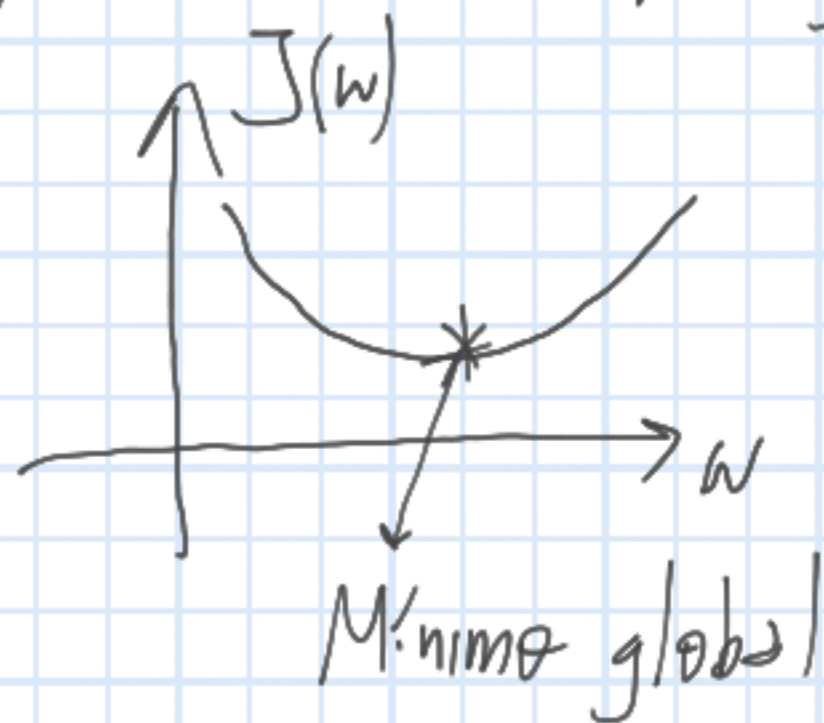


$$\hat{y} = f(x, b, w)$$

$$L = f(x_1, x_2, \dots, x_n, b, w) = \frac{1}{n} \sum_{i=1}^n \ell_i(x_i, b, w)$$

$$\text{MSE} (y_i - \hat{y})$$

* Algoritmos de optimización



$w_i \sim \text{Uniforme} \left(-\frac{1}{\sqrt{n}}, +\frac{1}{\sqrt{n}} \right)$ → cantidad de datos de entrada de la neurona

* Mini-Batch SGD

SGD con decaimiento lineal/
momento de primer orden

AdaGrad

RMS Prop

* → Adam (2014)

Variante $\bar{W} \leftarrow \bar{W} - \alpha \nabla J$
Derivada primer orden

* Newton

BFGS

L-BFGS

Derivada de
segundo orden
 H

* Mini-batch

for e in epochs:

for b in batches:

* Forward

$$* \bar{G} = \bar{\nabla}_{\bar{W}} \left(\frac{1}{b} \sum_{i=1}^{N\text{-BATCH}} \ell(\text{net}(\bar{x}_i, \bar{W}), y_i) \right)$$

$$* \bar{W} \leftarrow \bar{W} - \alpha \bar{G}$$

* First order momentum

for e in epochs:

for b in batches:

* Forward: $\hat{y} = \text{NET}(\bar{x})$

* \bar{G} (calcular el gradiente)

$$* \bar{V} \leftarrow \gamma \bar{V} + \alpha \bar{G}$$

$$* \bar{W} \leftarrow \bar{W} - \bar{V}$$

para todo el batch

$\bar{V} = 0 \rightarrow$ Vanilla Mini-Batch

$\bar{V} > 0$

¿Qué aprende? $\rightarrow \bar{W}$

¿Quiénes son HP? $\rightarrow \alpha, N\text{-BATCH}, N\text{-EPOCH}$

HP $\rightarrow \alpha, n\text{-epoch}, n\text{-batch}$
 \bar{V}

* AdaGrad

for e in epochs:

for b in batches:

* Forward

* \bar{G}

* $\bar{r} \leftarrow \bar{r} + \bar{G} \odot \bar{G}$

* $\bar{\Delta} \leftarrow \frac{-\alpha}{\sqrt{\bar{r}}} \odot \bar{G}$

* $\bar{w} \leftarrow \bar{w} + \bar{\Delta} = \begin{bmatrix} w_1 \\ \vdots \\ w_L \end{bmatrix} - \begin{bmatrix} \alpha_1 \frac{\partial J}{\partial w_1} \\ \vdots \\ \alpha_L \frac{\partial J}{\partial w_L} \end{bmatrix}$

element-wise
multiplication

$$\begin{bmatrix} \frac{\partial J}{\partial w_1} & \frac{\partial J}{\partial w_1} \\ \frac{\partial J}{\partial w_2} & \frac{\partial J}{\partial w_2} \\ \vdots & \vdots \end{bmatrix}$$

* Adam (2014)

for e in epochs:

for b in batches:

* Forward

* \bar{G}

* $\bar{V} \leftarrow \rho_1 \bar{V} + (1 - \rho_1) \bar{G}$

* $\bar{r} \leftarrow \rho_2 \bar{r} + (1 - \rho_2) \bar{G} \odot \bar{G}$

* $\bar{\Delta} \leftarrow -\alpha / \sqrt{\bar{r}} \cdot \bar{V}$

* $\bar{W} \leftarrow \bar{W} + \bar{\Delta}$