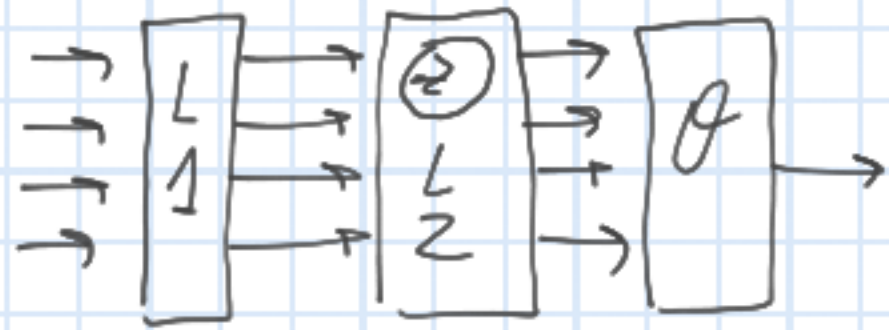


Funciones de activación

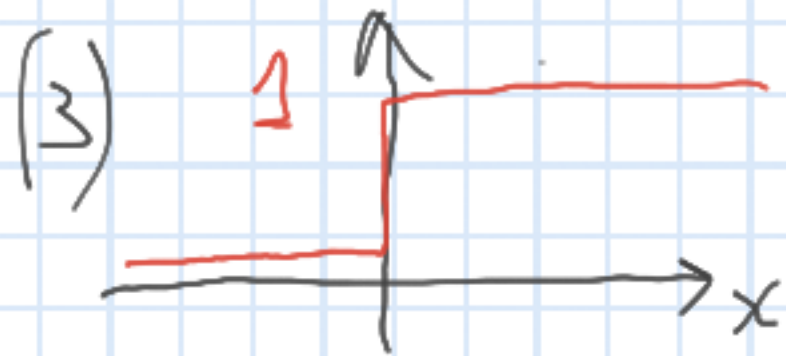
(1) No usar func. act.?



¿Qué problema tiene?

Esto es un modelo lineal

(2) Agregar func. act. / $\hat{y} = f(\bar{x})$ pueda ser cualquier f aprendida por la red



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

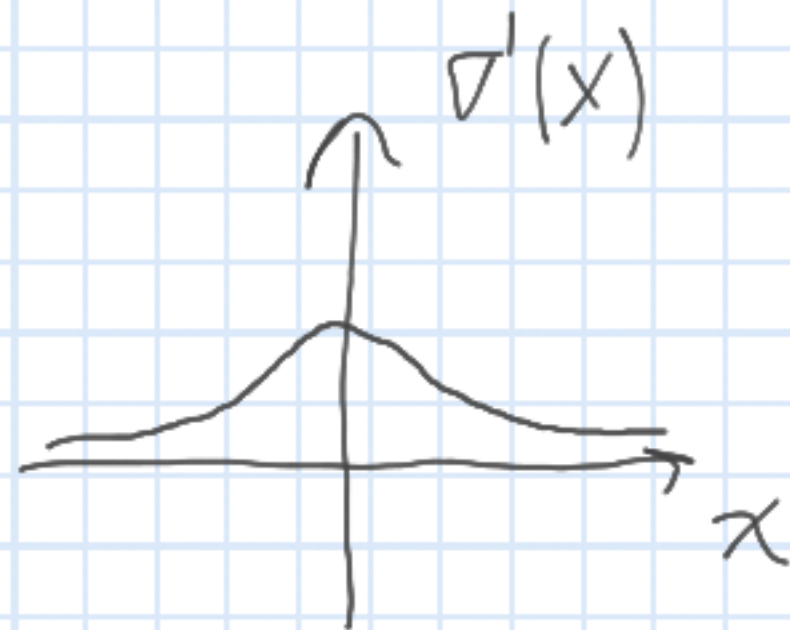
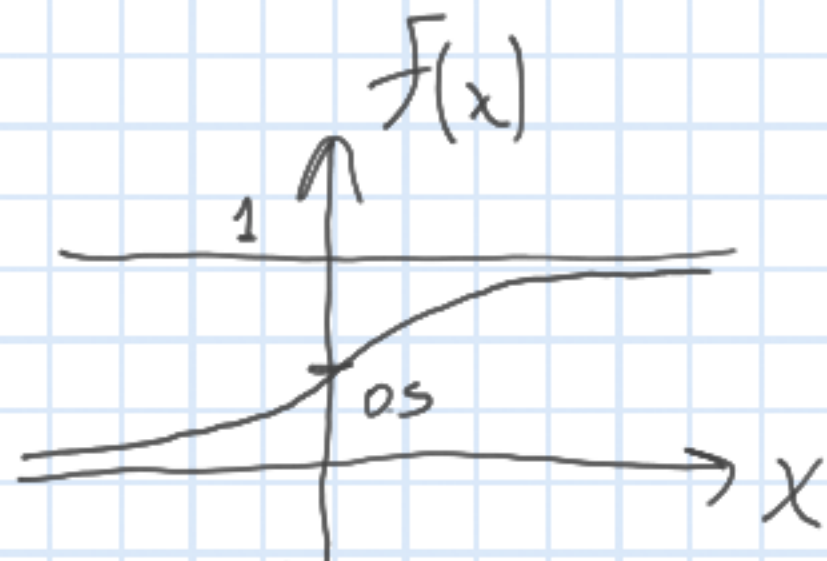
¿Qué problema tiene esta $f(x)$?

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

val = 0
cero!

(3) Sigmoid

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



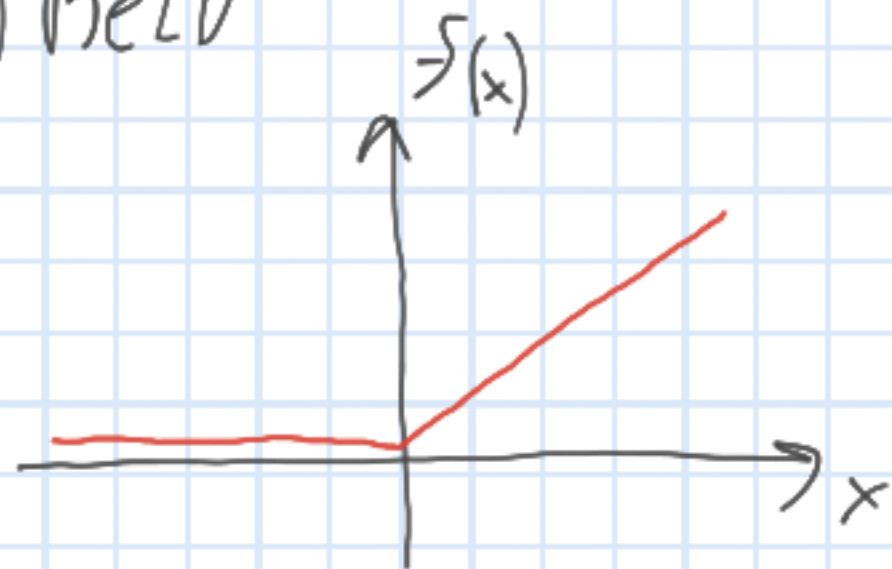
$$\lim_{x \rightarrow +\infty} \sigma'(x) = 0$$

$$\lim_{x \rightarrow -\infty} \sigma'(x) = 0$$

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

Problem: Vanishing gradients

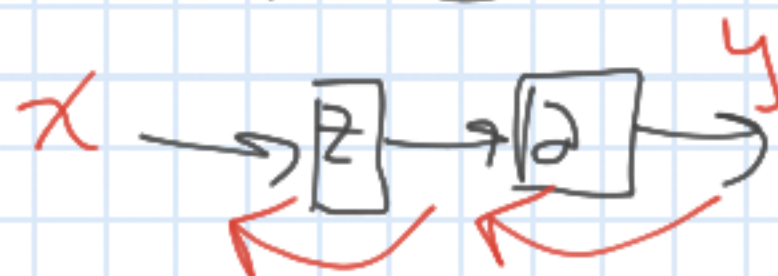
(4) ReLU



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

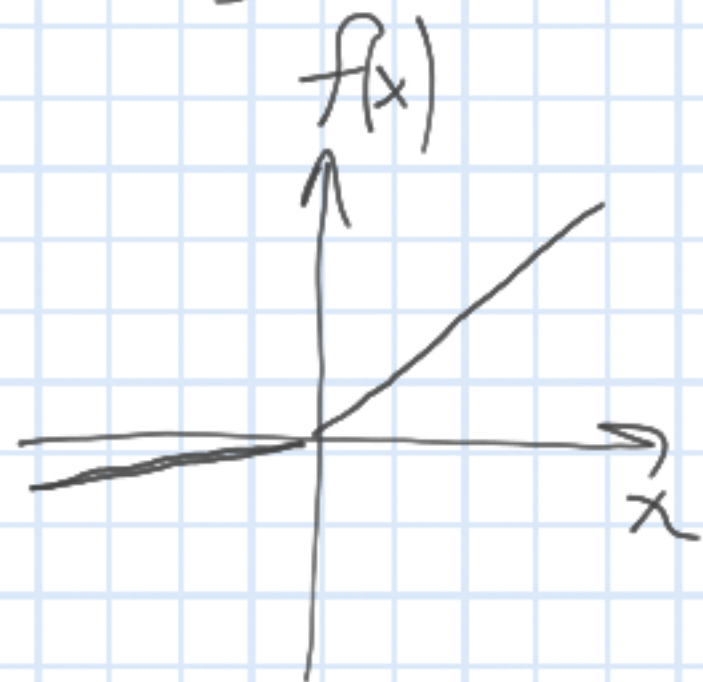
$$f'(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$$

$$z = w_1 x_1 + \dots + w_n x_n + b$$



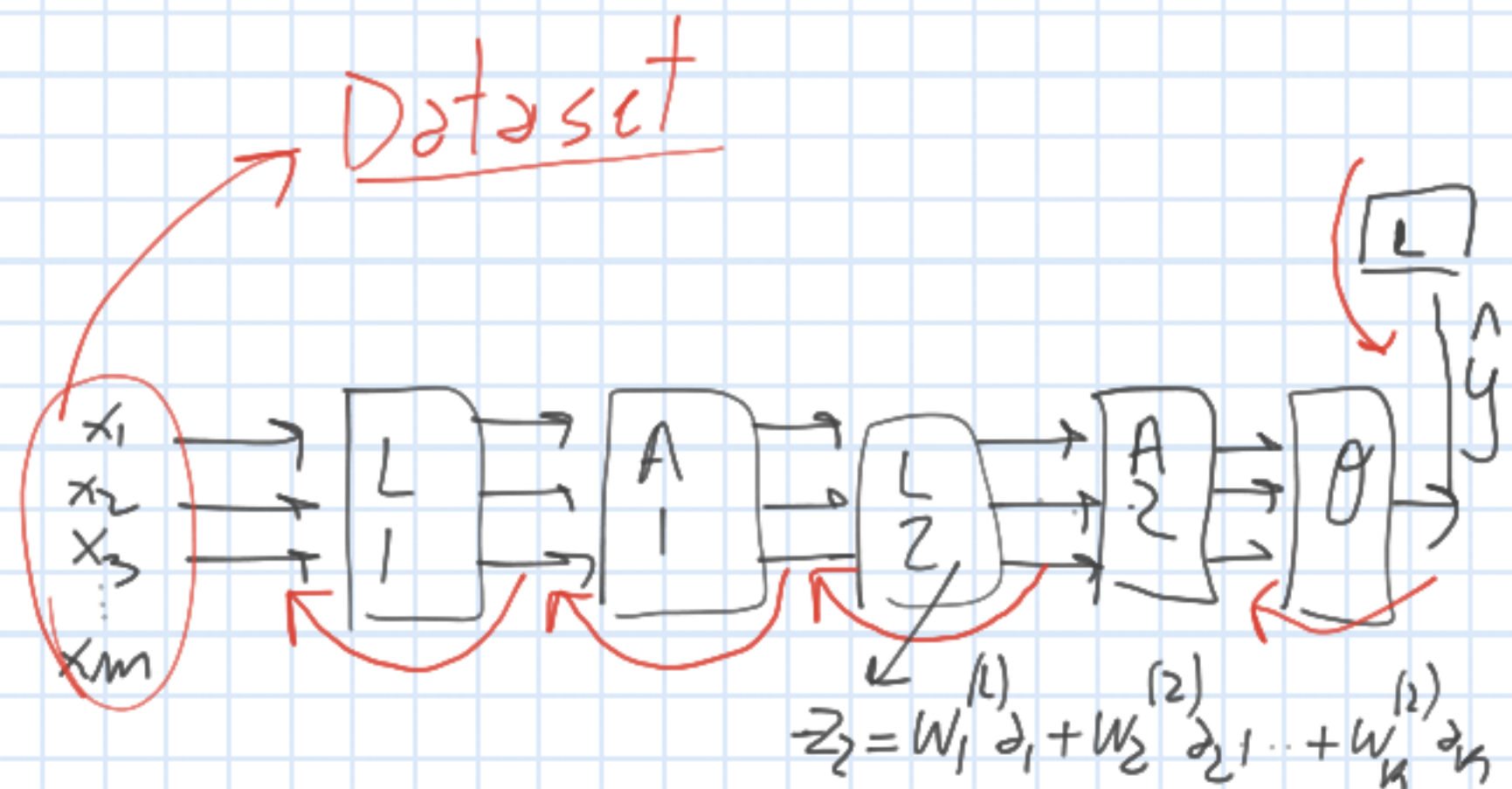
$$\frac{dy}{dx} = \frac{dy}{dz} \frac{dz}{dx}$$

(5) Leaky Relu



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

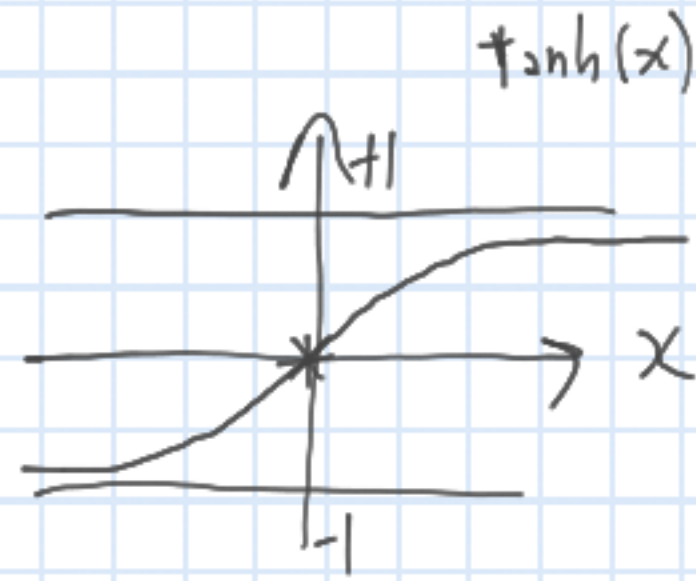
$$f'(x) = \begin{cases} 1 & x > 0 \\ 0.01 & x \leq 0 \end{cases}$$



$$\frac{dL}{dw_1^{(2)}} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{da_2} \frac{da_2}{dz_2} \frac{dz_2}{dw_1^{(2)}}$$

(3 con mejoras)

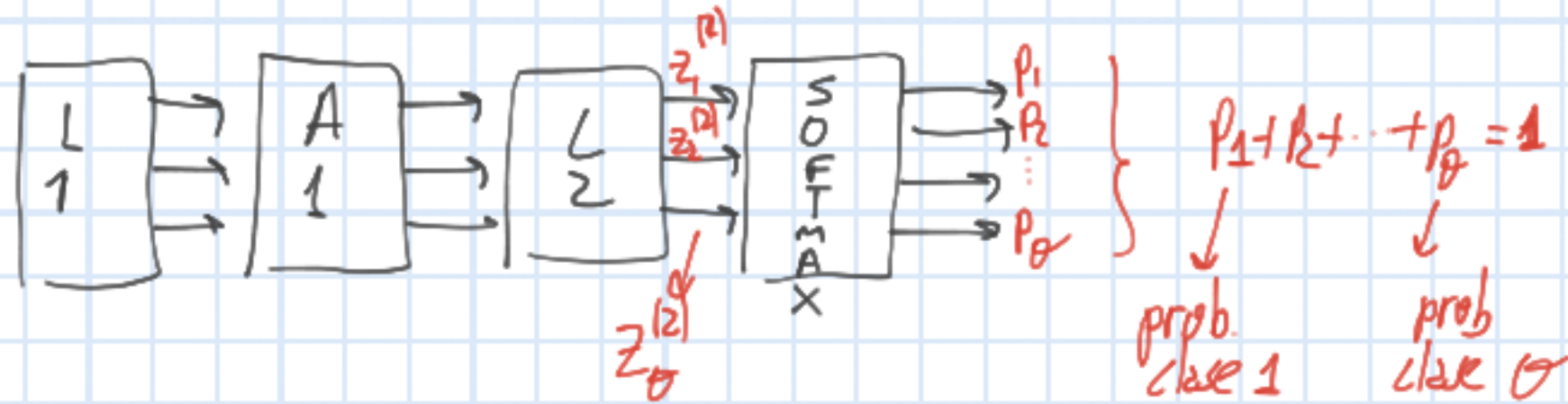
$$\tanh(x) = 2\sigma(x) - 1$$



Funciones de salida

Regresión ϕ
Clasif Bin σ

Clasif multi clase \rightarrow SOFTMAX

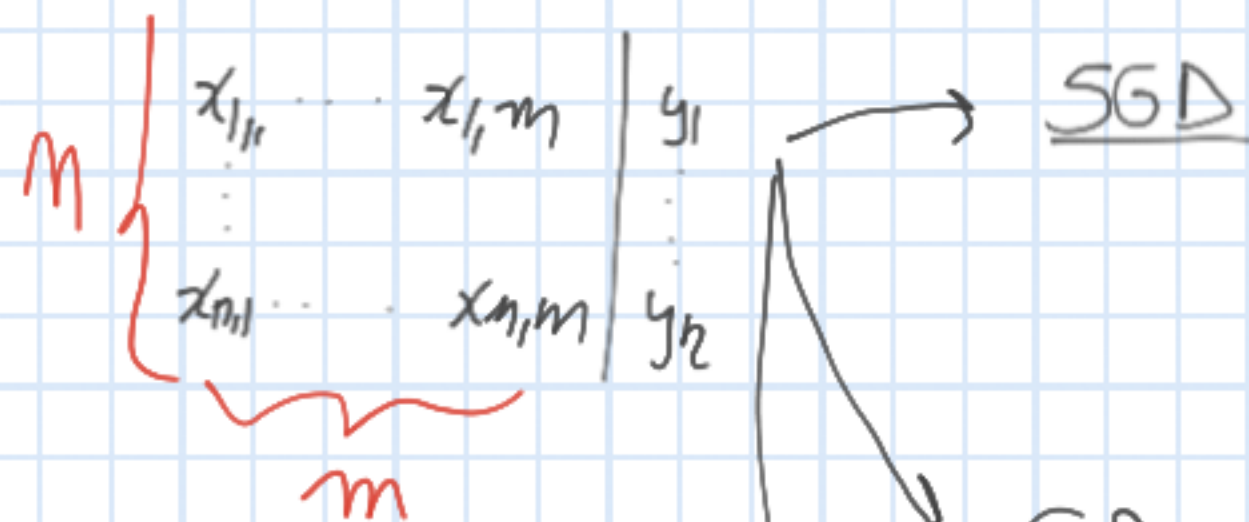


$$p_1 = \frac{e^{z_1}}{\sum_{i=1} e^{z_i}}$$

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

Loss Functions

(1) Binary Cross Entropy



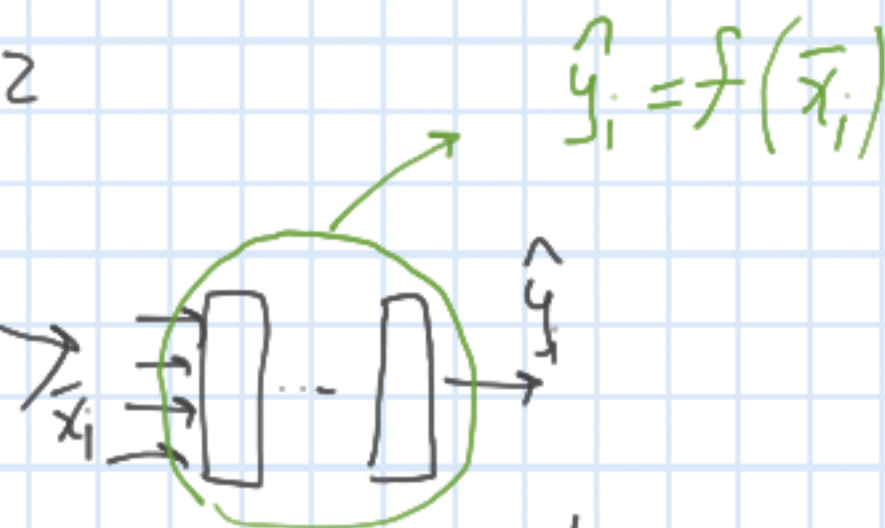
GD

ejemplo

$$* \text{Loss}(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2 = (y_i - f(\bar{x}_i))^2$$

↳ Computo gradiente
↳ Actualizar pesos

Representa la red neuronal



$$* \text{Loss}(y_1, y_2, \dots, y_n, \hat{y}_1, \dots, \hat{y}_n) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \rightarrow \begin{bmatrix} (y_1 - \hat{y}_1)^2 \\ \vdots \\ (y_n - \hat{y}_n)^2 \end{bmatrix}$$

↳ Computo gradi
↳ Actualizar pesos

$$\nabla(L) = \nabla_{\bar{w}} \left(\frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2 \right) = \frac{1}{n} \sum_{i=1}^n \nabla_{\bar{w}} [(y - \hat{y})^2]$$

	GD	Mini	SGD
Memoria	+++	++	+
Velocidad	+	++	+++

→ m-epochs * n

Mini-Batch → $b < n$

$$L = \frac{1}{b} \sum_{i=1}^b (y_i - \hat{y}_i)^2$$

(1) BCE

Modelar la salida en función de la entrada

$$P(\hat{Y}_i = 1 | \bar{X}_i = \bar{x}_i) = p, \quad p = \sigma(f(\bar{x})) \rightarrow \text{red neuronal}$$

$$Y_i | \bar{X}_i = \bar{x}_i \sim \text{Bernoulli}(p) \begin{matrix} \nearrow 1 \\ \searrow 0 \end{matrix}$$

$$P(Y_1 | \bar{X}_1 = \bar{x}_1 \cap Y_2 | \bar{X}_2 = \bar{x}_2 \cap \dots \cap Y_m | \bar{X}_m = \bar{x}_m) = \prod_{i=1}^m P(Y_i | \bar{X}_i = \bar{x}_i)$$

encontrar los \bar{w} que maximicen la observación

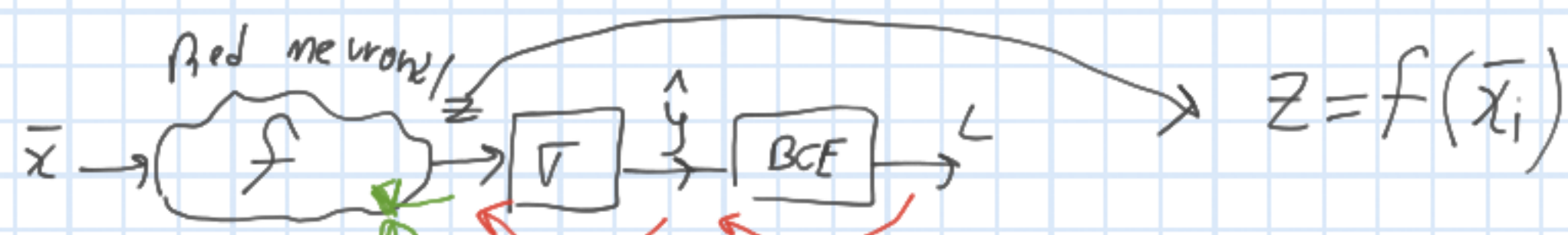
$$\arg \min_{\bar{w}} - \prod_{i=1}^m P(Y_i | \bar{X}_i = \bar{x}_i) \quad Y_i = y_i | \bar{X}_i = \bar{x}_i$$

$$\arg \min_{\bar{w}} - \frac{1}{n} \log \left(\prod_{i=1}^m P(Y_i | \bar{X}_i = \bar{x}_i) \right)$$

$$\arg \min_{\bar{w}} - \frac{1}{n} \sum_{i=1}^m \log (P(Y_i = y_i | \bar{X}_i = \bar{x}_i))$$

$$\begin{aligned} \arg \min_{\bar{w}} - \frac{1}{n} \sum_{i=1}^m \log (p_i^{y_i} (1-p_i)^{1-y_i}) \\ \arg \min_{\bar{w}} - \frac{1}{n} \sum_{i=1}^m \left[\log(p_i^{y_i}) + \log((1-p_i)^{1-y_i}) \right] \\ \arg \min_{\bar{w}} - \frac{1}{n} \sum_{i=1}^m \left[\underbrace{y_i \log(p_i) + (1-y_i) \log(1-p_i)}_{\text{BCE}} \right] \end{aligned}$$

$$\arg \min_{\bar{w}} - \frac{1}{n} \sum_{i=1}^m \left[y_i \log(\hat{y}_i) + \dots \right]$$



Derivation of the loss gradient with respect to the predicted output \hat{y} :

$$\frac{dL}{d\hat{y}} = - \left[y \log(\hat{y}) + (1-y) \log(1-\hat{y}) \right]$$

$$= -\frac{y}{\hat{y}} + \frac{(1-y)}{(1-\hat{y})}$$

The chain rule is applied as follows:

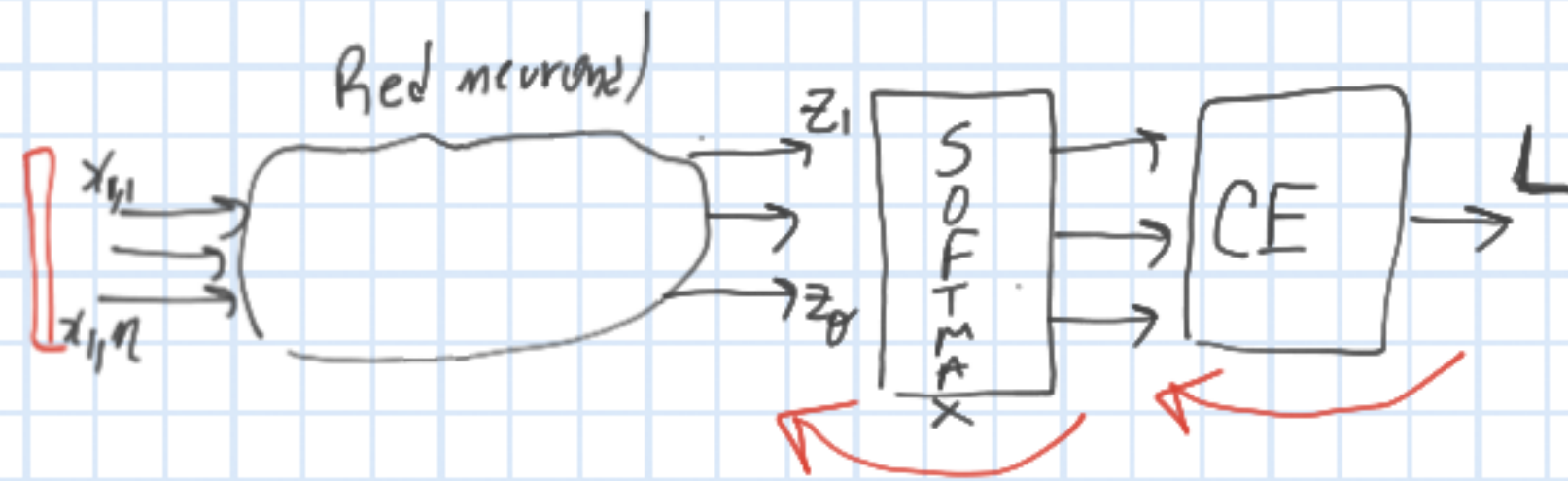
$$\frac{dL}{d\hat{y}} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dz} \cdot \frac{dz}{d\bar{x}}$$

Where the intermediate derivatives are:

$$\frac{d\hat{y}}{dz} = \sigma(z)(1-\sigma(z))$$

$$\frac{dL}{d\hat{y}} = -\left[y \log(\hat{y}) + (1-y) \log(1-\hat{y}) \right]$$

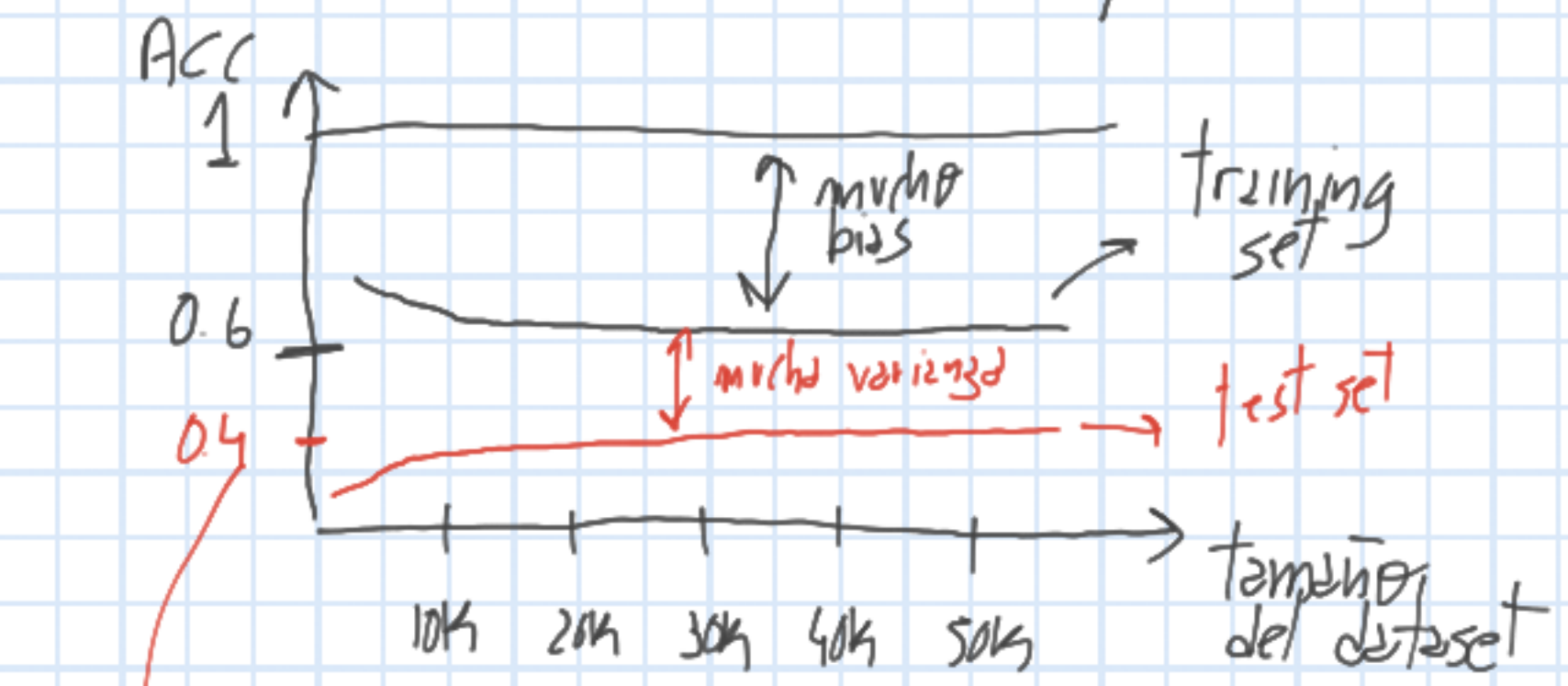
(2) Cross Entropy \rightarrow Class. multiclass



TODO

Regularización

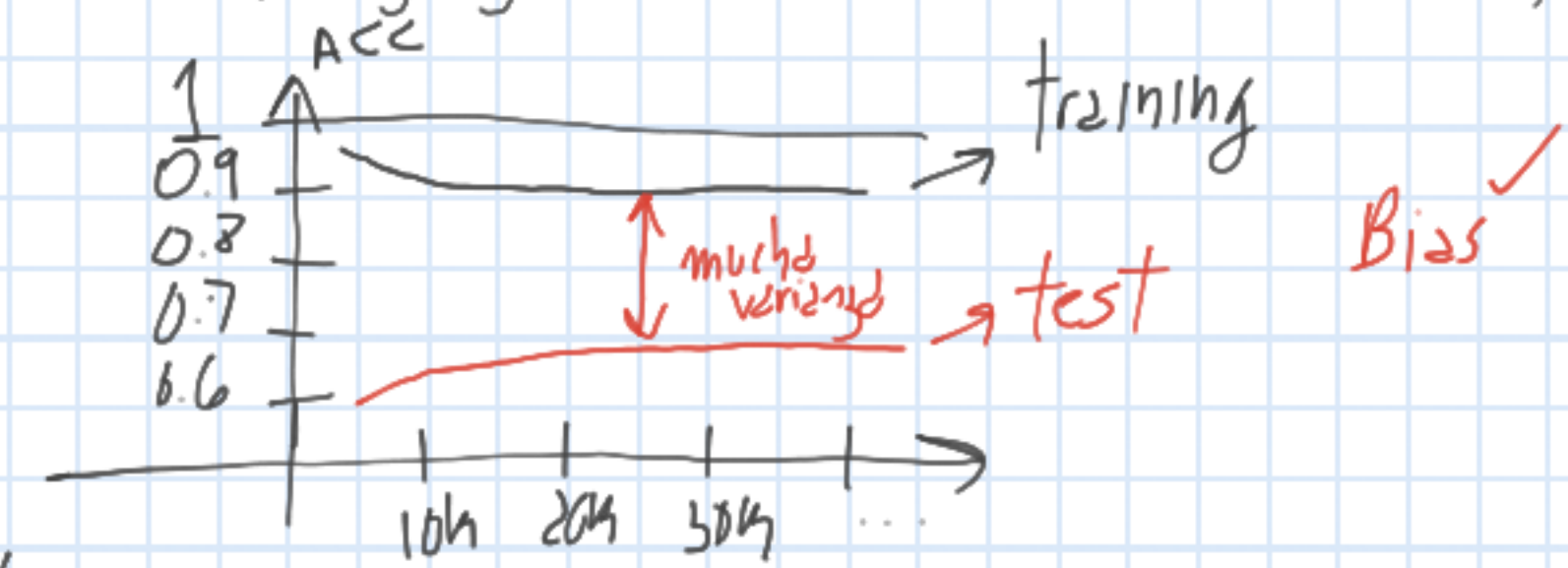
Bias vs Varianza
Underfit. vs Overfit.



Learning curve

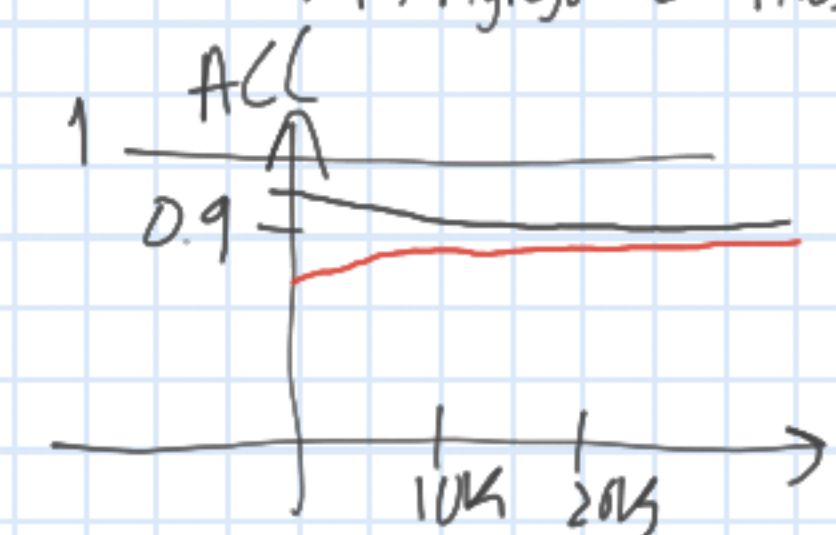
¿Cómo disminuir el error de bias?

- (1) Agregar complejidad al model
- (2) Agregar más features (+ columnas)



¿Cómo bajar el error por varianza?

- (1) Regularización
- (2) Agregar más filas (data)



* Regularización $\begin{cases} L2 \\ L1 \\ Dropout \end{cases}$

L2

$J(\bar{w}, y, \bar{x}) \rightarrow$ es un escalar \rightarrow sin regularizar

$$\tilde{J}(\bar{w}, y, \bar{x}, \lambda) = J(\bar{w}, y, \bar{x}) + \frac{\lambda}{2} \bar{w}^T \bar{w}$$

$$\bar{w}^T \bar{w} = [w_1 \dots w_n] \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} = w_1^2 + w_2^2 + \dots + w_n^2$$
$$= \|\bar{w}\|_2^2$$

todos los
pesos de
mi red neuronal

(1) Efecto sobre update a los pesos

$$\nabla_{\bar{w}}(\tilde{J}) = \nabla_{\bar{w}}\left(J + \frac{\lambda}{2} \bar{w}^T \bar{w}\right) = \nabla_{\bar{w}}(J) + \nabla_{\bar{w}}\left(\frac{\lambda}{2} \bar{w}^T \bar{w}\right) = \nabla_{\bar{w}}(J) + \lambda \bar{w}$$

Actualización del parámetro conocido

$$\bar{w} \leftarrow \bar{w} - \alpha \nabla(\tilde{J})$$

$$\leftarrow \bar{w} - \alpha (\nabla(J) + \lambda \bar{w})$$

$$\leftarrow \bar{w} - \alpha \lambda \bar{w} - \nabla(J) \alpha$$

$$\leftarrow (1 - \lambda \alpha) \bar{w} - \alpha \nabla(J)$$

(2) Efecto sobre el modelo

$$\hat{y} = \bar{w}^T \bar{x} \xrightarrow{\text{MSE}} \bar{w} = (X^T X)^{-1} X^T y$$

$$\text{MSE} + L2 \longrightarrow \bar{w} = \left(\underbrace{X^T X}_{\substack{\text{ETB} \\ m \times n}} + \underbrace{\lambda I}_{\substack{\text{ETB} \\ n \times n}} \right) X^T y$$

$\text{ETB } (m \times m)$

L1

$$\tilde{J} = J + \lambda \|\bar{w}\|_1 \quad \|\bar{w}\|_1 = |w_1| + |w_2| + \dots + |w_n|$$

$$\nabla(\tilde{J}) = \nabla(J) + \lambda \operatorname{sign}(\bar{w})$$



$$\bar{w} \leftarrow \bar{w} - \alpha \nabla(\tilde{J})$$

$$\bar{w} \leftarrow \bar{w} - \alpha (\nabla(J) + \lambda \operatorname{sign}(\bar{w}))$$

$$\bar{w} \leftarrow \bar{w} - \alpha \lambda \operatorname{sign}(\bar{w}) - \alpha \nabla(J)$$

↗ bajar complejidad
modelo