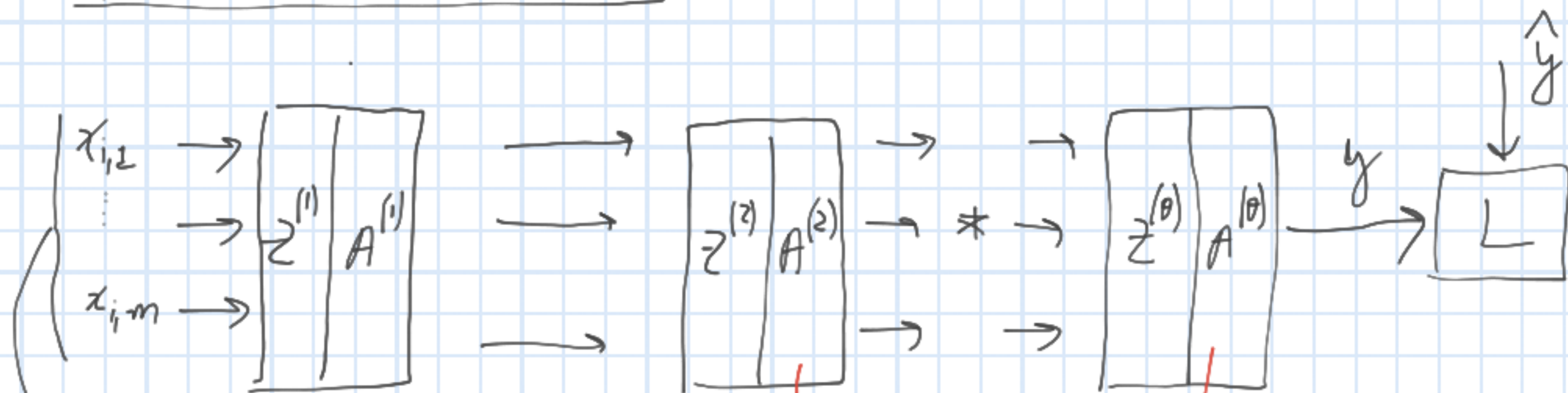


Clase 2:

* Funciones de activación



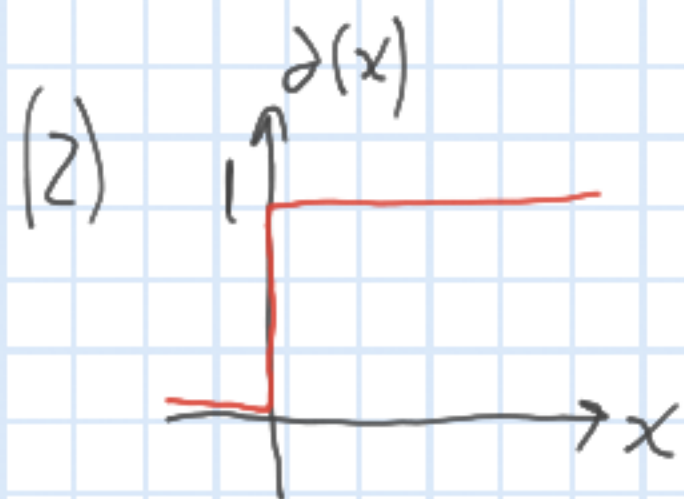
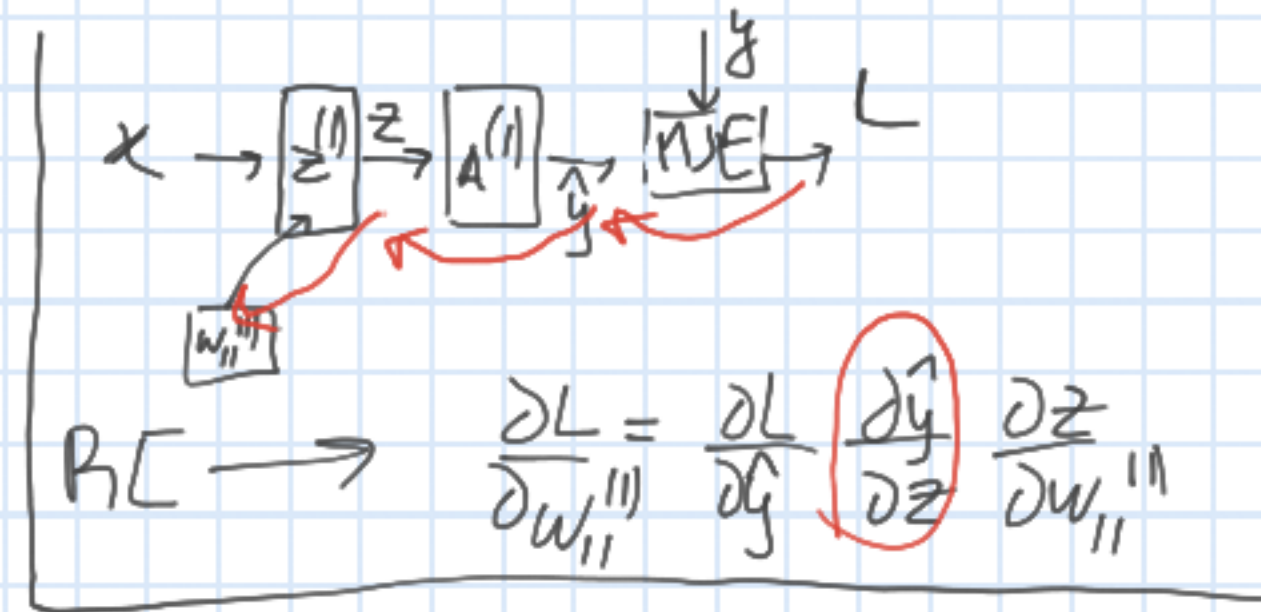
$X \in \mathbb{R}^{m \times 1}$ (#batch size $\rightarrow 1$)
 $X \in \mathbb{R}^{m \times b}$ (#batch size $\rightarrow b$)

Sigmoid
ReLU
tanh

SOFTMAX \rightarrow multi-class
Sigmoid \rightarrow clasificación binaria
Ninguna \rightarrow regresión

(1) No usar función de activación

$$x \rightarrow z^{(1)} \rightarrow z^{(2)} \rightarrow \hat{y} = f(x) = z^{(2)}(z^{(1)}(x))$$

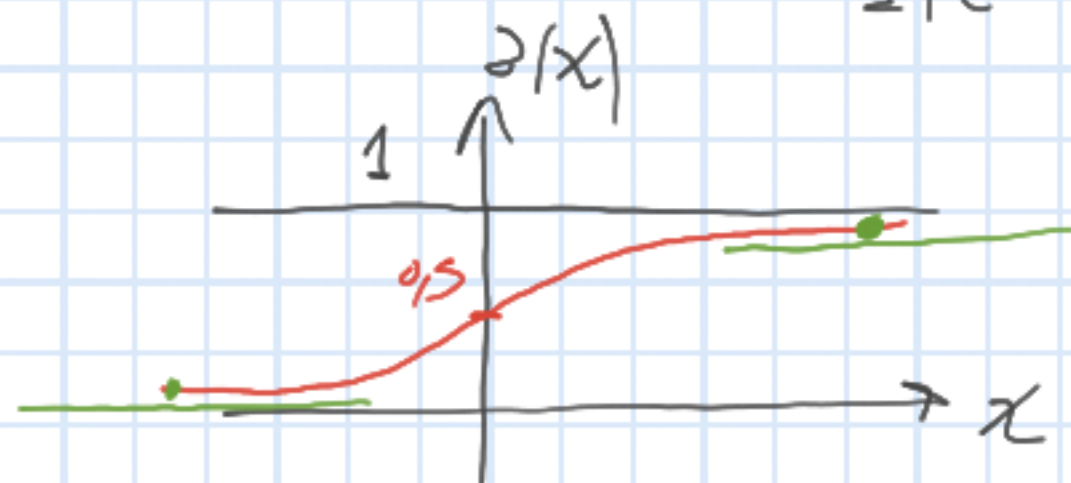


$$\sigma(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$$

$$\sigma'(x) = \begin{cases} 0 & x > 0 \\ 0 & x < 0 \\ ? & x = 0 \end{cases}$$

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

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

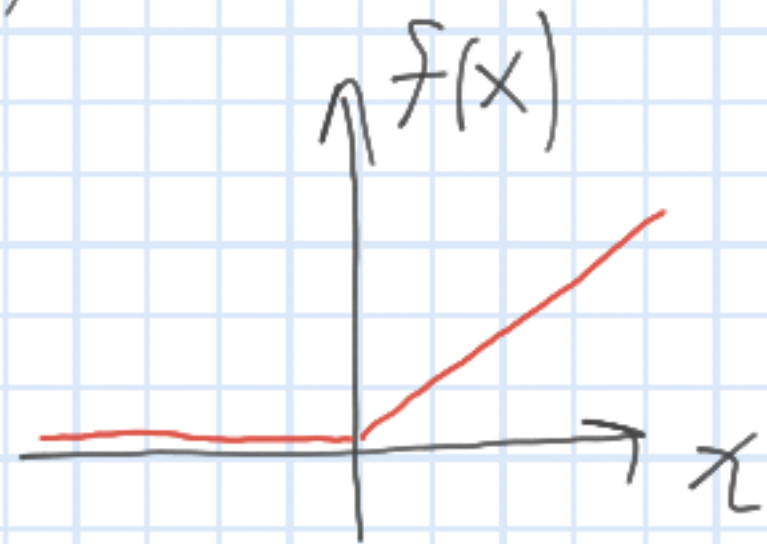


$$\sigma'(x) = \frac{\partial \tau(x)}{\partial x} = \tau(x)(1-\tau(x))$$

↳ Optimización lenta
Vanishing gradient



(4) ReLU



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

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

↳ Para $x < 0$, la derivada es zero

(5) Leaky ReLU



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

↳ $\alpha = 0.01$

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

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



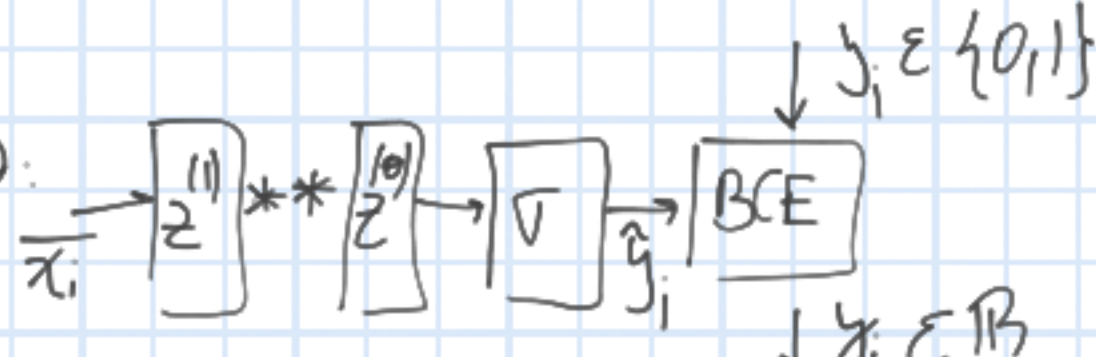
* Funciones salida + Loss Function

$$\left\{ \begin{array}{c} x_{1,1} \dots x_{1,m} \\ \vdots \\ x_{n,1} \dots x_{n,m} \end{array} \middle| \begin{array}{c} y_1 \\ \vdots \\ y_n \end{array} \right\} \rightarrow f_{y_i/\bar{x}_i}(y_i/\bar{x}_i)$$

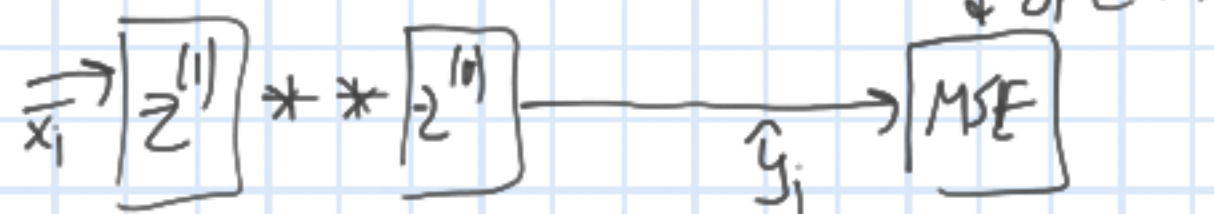
$$L(\bar{x}_1, \dots, \bar{x}_n, y_1, \dots, y_n) = \mathbb{P}(Y_1=y_1, Y_2=y_2, \dots, Y_n=y_n \mid \bar{X}_1=\bar{x}_1, \dots, \bar{X}_n=\bar{x}_n)$$

$$= \prod_{i=1}^n \mathbb{P}(Y_i=y_i \mid X_i=x_i) \rightarrow \prod_{i=1}^n f_{y_i/\bar{x}_i}(y_i/\bar{x}_i)$$

* Clasificación Binaria:



* Regresión:



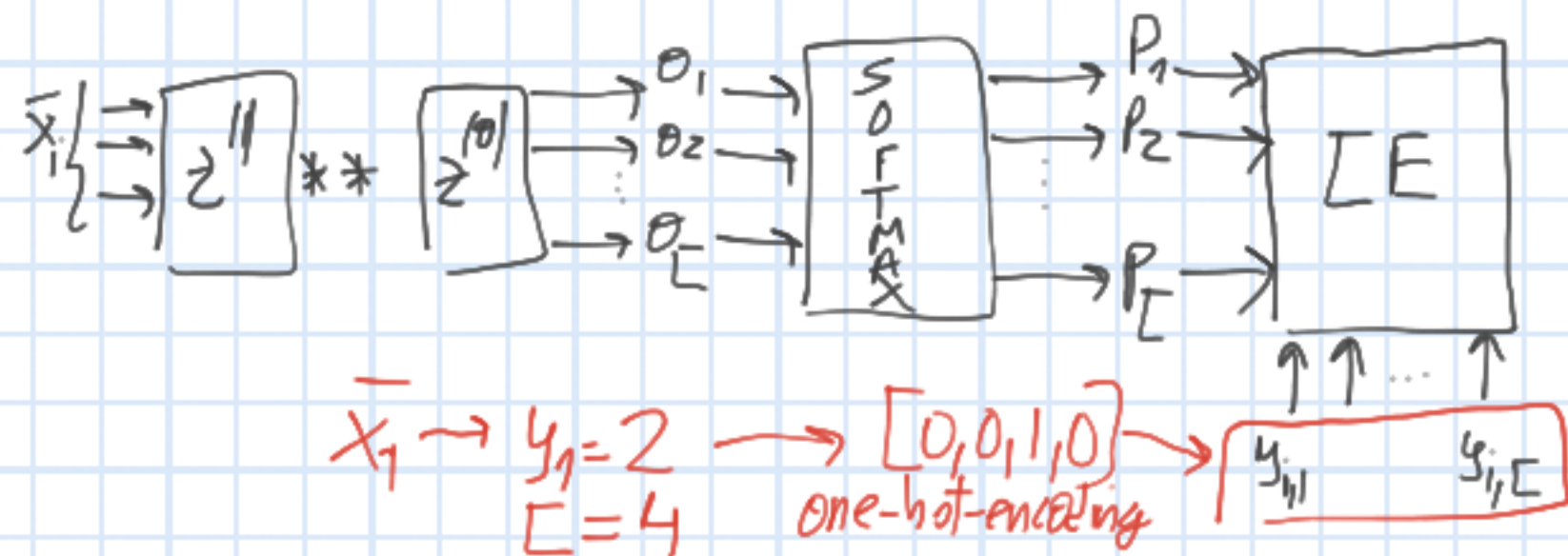
(1) Clasificación binaria $\rightarrow Y_i/\bar{X}_i \sim \text{Bernoulli}(p_i)$
 $p_i = \bar{w}^T \bar{x}_i$

(2) Regresión $\rightarrow \text{Loss} \rightarrow \text{BCE}$

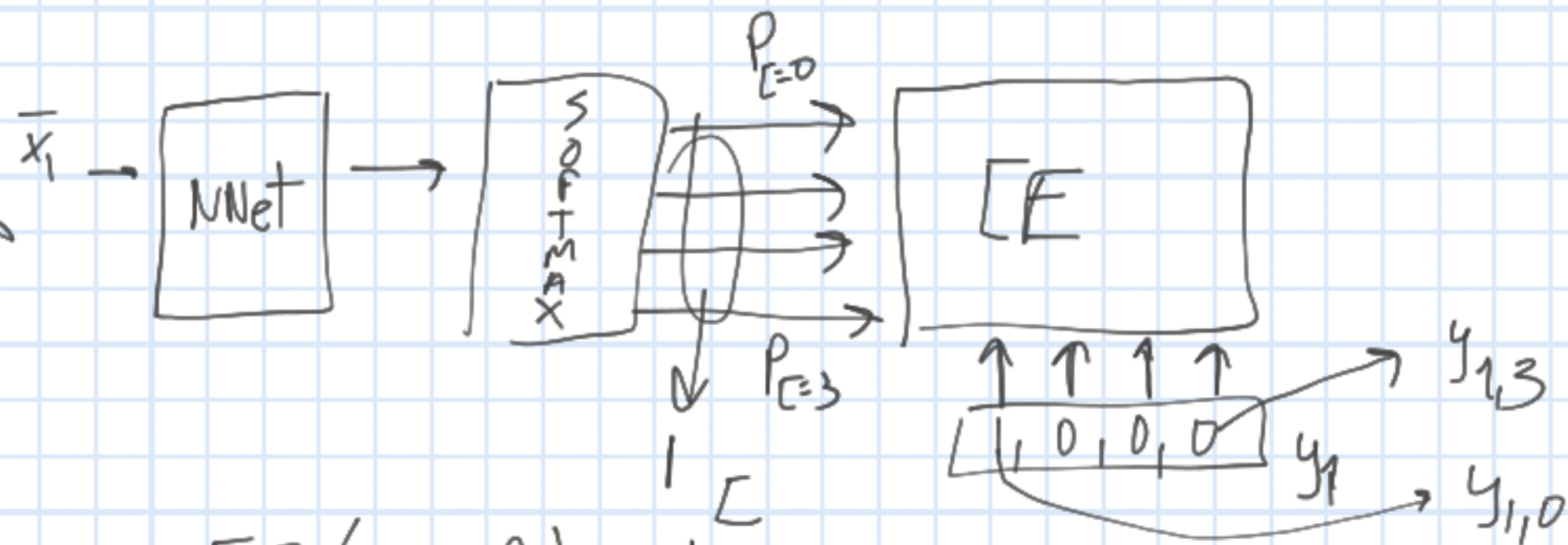
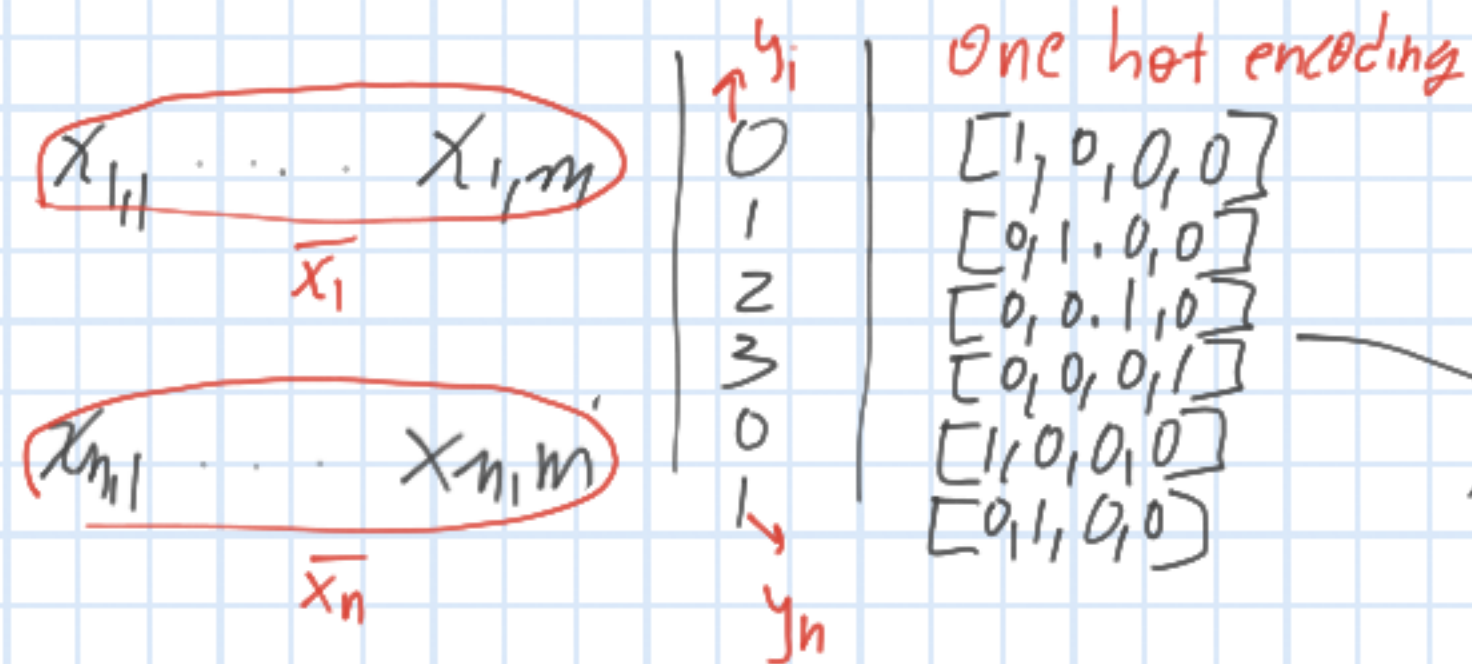
$\rightarrow Y_i/\bar{X}_i \sim \text{Normal}(\mu_i, 1)$, $\mu_i = \bar{w}^T \bar{x}_i$

$\rightarrow \text{Loss} \rightarrow \text{MSE}$

* Clas. multiclase



$$p_i = \frac{e^{\theta_i}}{\sum_{j=1}^L e^{\theta_j}} \quad \left| \quad \sum_{i=1}^L p_i = 1 \right.$$

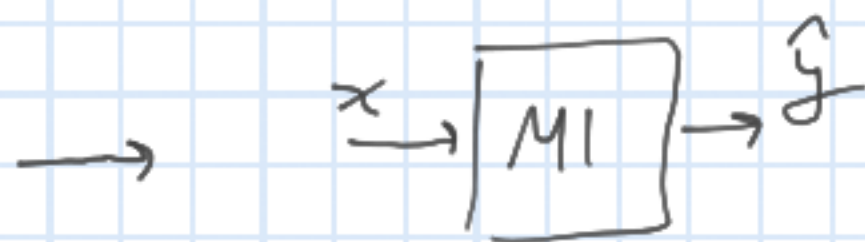
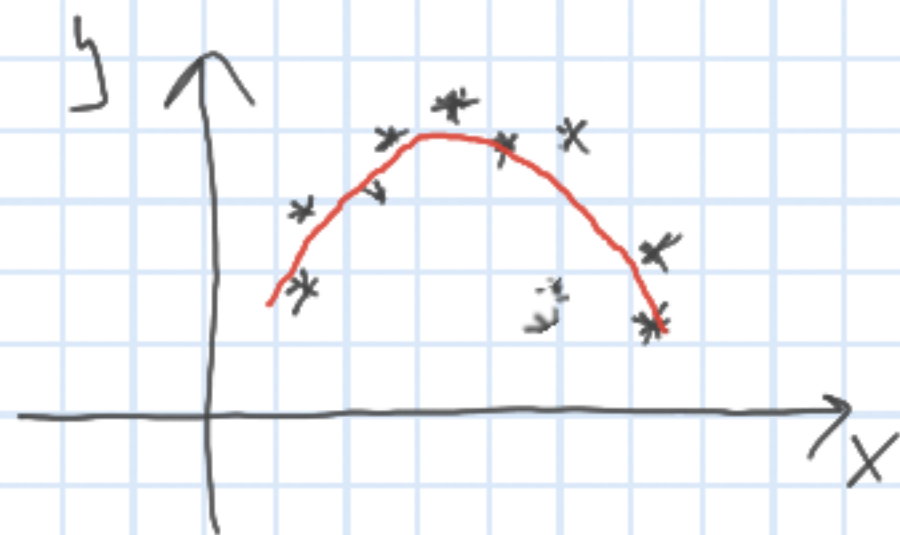


$$[E](y_i, \hat{y}_i) = \sum_{j=0}^L y_{i,j} \log(P_{i,j})$$

$$= \sum_{j=0}^L y_{i,j} \log \left(\frac{e^{\theta_{i,j}}}{[TE_i]} \right)$$

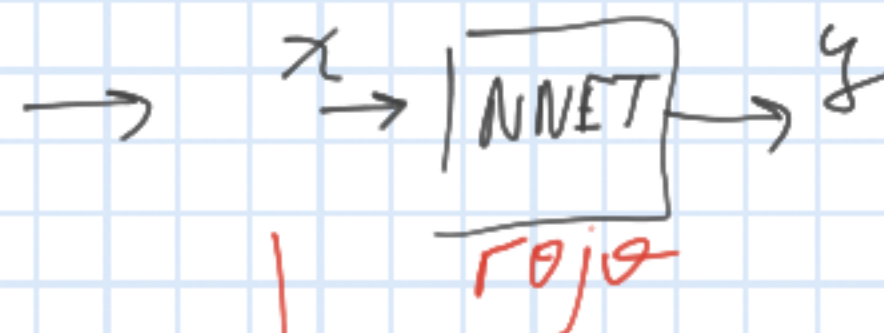
$$= \sum_{j=0}^L y_{i,j} \left[\log(e^{\theta_{i,j}}) - \log([TE_i]) \right]$$

$$= \sum_{j=0}^L y_{i,j} \left[\theta_{i,j} - \log([TE_i]) \right]$$

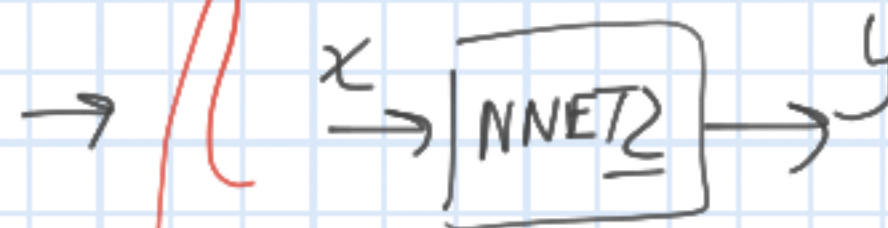
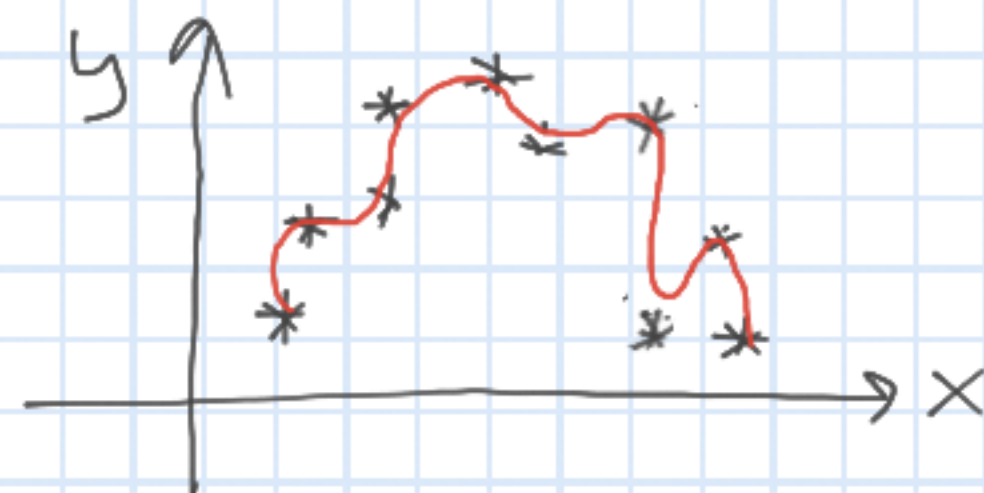


$$\hat{y}_x = w_1 x^2 + w_2 x + w_3 \rightarrow \text{params } 3! \text{ (rojo)}$$

Ejercicio Clase 2 \rightarrow Usando SGD, encontrar w_1, w_2, w_3



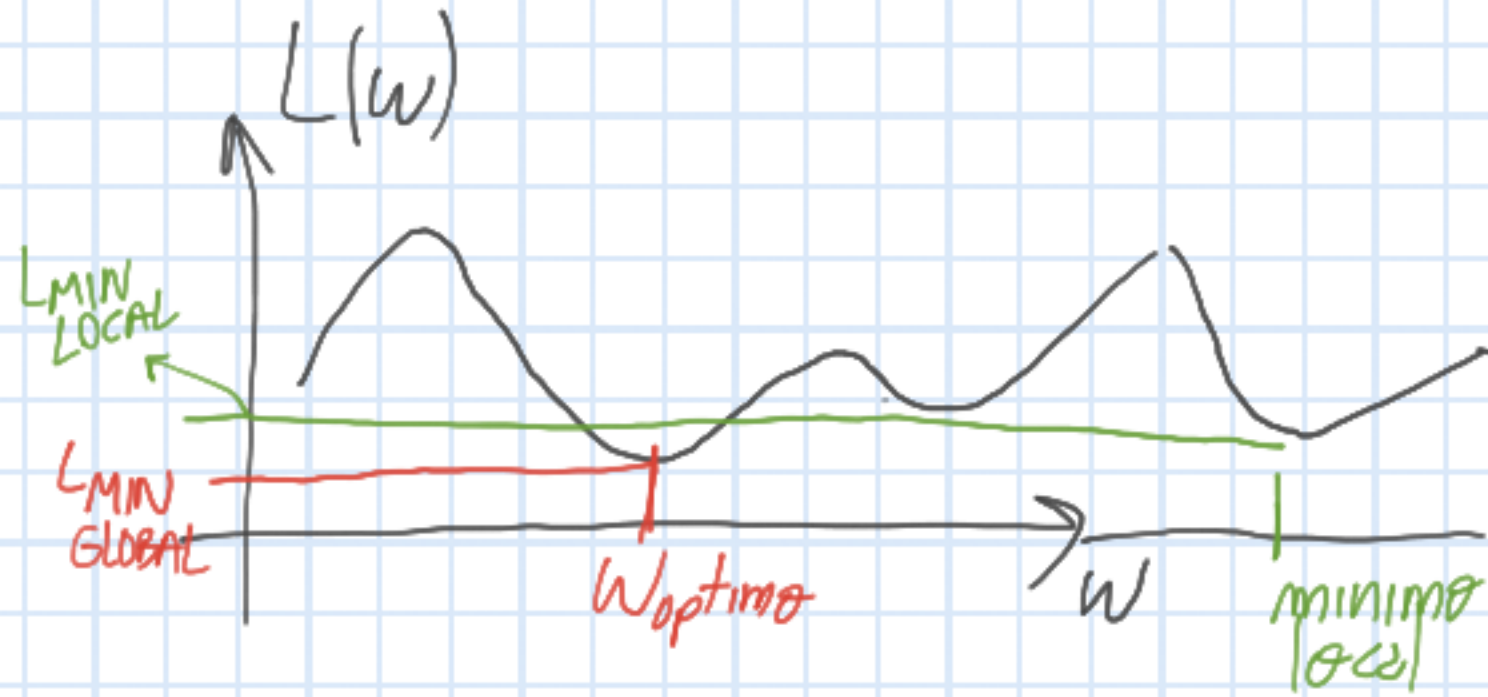
$\rightarrow \text{NNET} \rightarrow \text{FIT CORRECTO}$



$\rightarrow \text{OVERFITTING}$

\downarrow Clase 3 \rightarrow PyTorch

* Optimización



Optimizar
con gradiente
 $\bar{\nabla}(L)$

Optimizadores
que usan derivadas
de orden superior $H(L)$

Optimización
mediante búsquedas → Genética
→ Hill Climb

*
*

Mini-BATCH (GD MB)
GD usando decaimiento
GD momento primer orden
AdaGrad, RMSProp
Adam (2014)

Newton
BFGS
L-BFGS

* Mini-Batch

for e in epochs: \rightarrow epoch=10

for b in batches \rightarrow 64 muestras
bs=8

* \bar{X}_b
* $\hat{Y}_b = \text{Forward}(\bar{X}_b)$ } * FORWARD

* Error = $L(Y_b, \hat{Y}_b)$ } * LOSS FUNCTION

* $\bar{G} = \bar{\nabla}_W \left(\frac{1}{\#bs} \sum_{i=1}^{\#bs} \ell(y_i, \hat{y}_i) \right)$ \rightarrow $bs=1 \rightarrow$ SGD
 $1 < bs < n \rightarrow$ Mini-Batch \rightarrow Backward
 $bs=n \rightarrow$ GD

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

1* ¿Cuántas veces actualiza \bar{W} en una epoch?

$$64 / bs = 64 / 8 = \underline{8}$$

2* ¿Cuántas veces actualiza \bar{W} en total?

$$8 * ne = 8 * 10 = \underline{80}$$

\rightarrow ¿HPs?
optimizer \rightarrow α
epochs
 \rightarrow #bs

* AdaGrad

for e in epochs:

for b in batches:

* FORWARD

* $\bar{G} \leftarrow \text{BACKWARD}$

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

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

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

* Adam (2014)

* for e in epochs

* for b in batch

* Forward

* Backpropagation

* \bar{G}

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

momento 1st order

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

momento 2nd order

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

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